Quantifying success in individual and team careers in intellectual domains

by Sandeep Chowdhary

Supervisor: Federico Battiston

A Dissertation Submitted in Partial Fulfillment of the Requirements for the Degree of Doctor of Philosophy in Network Science



Central European University Vienna, Austria

Date of submission: February, 2025

Sandeep Chowdhary: *Quantifying success in individual and team careers in intellectual domains*, © 2025 All rights reserved.

RESEARCHER DECLARATION

I, Sandeep Chowdhary, certify that I am the author of the work "Quantifying success in individual and team careers in intellectual domains". I certify that this is solely my own original work, other than where I have clearly indicated, in this declaration and in the thesis, the contributions of others. The thesis contains no materials accepted for any other degrees in any other institutions. The copyright of this work rests with its author. Quotation from it is permitted, provided that full acknowledgement is made. This work may not be reproduced without my prior written consent.

CEU eTD Collection

Abstract

Human progress is driven by actions of successful individuals. For this reason, quantifying patterns and habits of success in major human endeavors is a key challenge in computational social science. Yet, its quantification is often challenging due to its subjective nature, and the lack of fine-grained data about human performance. The recent availability of large scale datasets in a variety of social contexts- from Science to Businesses and Sports- presents an opportunity for success to be rigorously quantified. The first part of this thesis explores temporal patterns of success in individual careers. First, I unveil how social network and collaboration structure determines funding success of scientists in academia. Furthering our data-driven approaches to sports, I delve into the evolution of chess careers, tracking the career trajectories of successful chess players, and their temporal evolution. Beyond single individuals, I then move on characterizing success in teams. Focusing on publication data, I extract persistent collaborations and extend the investigation of scientific careers from single scientists to team trajectories. I investigate how persistent collaborations emergence and eventually dissolve, and which compositional factors make some teams more successful than others. Lastly, I explore the spread of ideas and innovations through social contagion modeling, highlighting the importance of group contagion and temporal persistence in shaping populationlevel spread of ideas. By quantifying patterns success across careers in different domains, this thesis contributes to the evolving discourse on successful careers and their correlates.

CEU eTD Collection

ACKNOWLEDGEMENTS

I would like to acknowledge several individuals who have significantly influenced this work and its completion. First, I must thank my advisor, Federico Battiston, for his continued support and confidence. His guidance has been invaluable, and his adaptive supervision throughout my PhD has been essential. I would like to extend my gratitude to my friends and colleagues from the DNDS PhD cohort, as well as to the seniors for their guidance during my PhD. Additionally, I am thankful to the faculty at CEU for the guidance and feedback they provided in seminars and one-on-one meetings throughout my PhD. In particular, I am grateful to Gerardo Iñiguez for initiating an off-beat and exciting collaboration involving mostly students from a course he taught, which ultimately became a published paper.

I would also like to thank my collaborator, Federico Musciotto, for sharing his skills in data wizardry.

The support of my family throughout this endeavor has been an unwavering source of strength, carrying me through many challenges I faced during this long journey.

CEU eTD Collection

CONTENTS

C	onten	ts	i
Li	st of :	Figures	v
1	Intr	oduction	1
2	Rol ing	e of collaboration network in securing individual scientific fund-	7
	2.1	Introduction	7
	2.2	Data	8
	2.3	Cross-continental collaborations of ERC and NSF awardees	8
	2.4	US collaborations before and during an ERC grant	11
	2.5	Pre-award mobility of ERC vs NSF awardees	13
	2.6	Wider collaboration patterns of ERC and NSF funded researchers	15
	2.7	Universality of the EU-US imbalance	16
	2.8	The case of physics	17
	2.9	Discussion	19
3	Quantifying performance in chess careers		23
	3.1	Introduction	23
	3.2	Methods	24
	3.3	Careers in chess	28
	3.4	Hot and cold streaks in chess careers	28
	3.5	Effect of player skill on likelihood of hot-streaks	32
	3.6	Specialization in opening play and skill level	35
	3.7	Performance and opening diversity	39
	3.8	Diversity vs Career stage	41
	3.9	Discussion	41

Team careers in science: formation, composition and success of persis-			
tent collaborations			
4.1	Introduction	45	
4.2	Methods	46	
4.3	Formation, productivity and dissolution of persistent scientific		
	collaborations.	50	
4.4	Composition.	53	
4.5	Team success.	60	
4.6	Discussion	66	
Model of social contagion in a population with evolving group struc-			
ture		71	
5.1	Introduction	71	
5.2	Model of contagion	73	
5.3	Social contagion on static and temporal simplicial complexes	75	
5.4	Contagion on temporally correlated higher-order networks	79	
5.5	Contagion on degree-heterogeneous temporal higher-order net-		
	works	81	
5.6	Discussion	84	
Con	clusion	87	
	Tean tentt 4.1 4.2 4.3 4.4 4.5 4.6 Mood ture 5.1 5.2 5.3 5.4 5.5 5.6 Con	Team careers in science: formation, composition and success of persistent collaborations 4.1 Introduction 4.2 Methods 4.3 Formation, productivity and dissolution of persistent scientific collaborations. 4.4 Composition. 4.5 Team success. 4.6 Discussion 4.6 Discussion 5.1 Introduction 5.2 Model of contagion in a population with evolving group structure 5.1 Introduction 5.3 Social contagion on static and temporal simplicial complexes 5.4 Contagion on temporally correlated higher-order networks 5.5 Contagion on degree-heterogeneous temporal higher-order networks 5.6 Discussion	

LIST OF FIGURES

2.1	Asymmetric collaborations of European and American funded research.	9
2.2	Collaboration imbalances between EU based scientists and US based scientists.	9
2.3	Differences in collaboration patterns before and after winning	10
	the grant.	10
2.4	Number of collaborators before and during an ERC grant.	11
2.5	Comparison between ERC winners and top NSF awardees.	12
2.6	Asymmetric mobility of European and American funded re-	
	searchers.	13
2.7	Mobility of scientists across countries before grant award	14
2.8	Mobility differences before grant award: EU as a single state.	14
2.9	Ratio between percentages of ERC and NSF collaborations	
	with EU countries.	15
2.10	Ratio between percentages of ERC and NSF collaborations	
	with non-EU countries.	16
2.11	Discipline-wise quantification of asymmetrical cross- continental collaborations patterns for ERC and NSF awardees.	
	1	17
2.12	Collaboration imbalance across scientific domains.	18
2.13	Co-authorship and mobility differences between funded	
	physicists in the EU and the US.	19
2.14	Countrywise collaborations of ERC and NSF winning physicists.	20
3.1	Distribution of player ratings in data.	27
3.2	Distribution of the number of games per player.	28
3.3	Visualization of the career of the Grandmaster (GM) Magnus	
	Carlsen.	29
3.4	Hot and cold streaks in chess careers.	30
3.5	Autocorrelation test for hot-treaks.	30
3.6	Role of rating advantage over opponent in hot-streaks.	31

3.7	Rating disadvantage in the hot-streak breaking game.	32
3.8	Role of time-difference between successive games in hot-streaks	33
3.9	Distribution of number of repeated matches in the dataset.	33
3.10	Distribution of balance in outcomes in player match-ups for	
	the 4 skill categories.	34
3.11	Hot-streaks and player skill.	34
3.12	Diversity and specialization in the first move of the game.	35
3.13	Diversity in the black's response to white's first move.	36
3.14	Top 9 favorite openings among all players on lichess.	37
3.15	Usage of top-2 favourite openings of a player.	38
3.16	Diversity and specialization in the opening sequence of moves	
	is governed by skill level.	38
3.17	Opening switches in a player's career properly normalised	
	with the null model.	39
3.18	Winrate of top 3 openings of a player against opening frequency.	40
3.19	Sub-optimal opening encounters.	40
3.20	Diversity in opening with career stage.	42
4.1	Identification of persistent core members in a team's career	50
4.2	Distribution of core sizes.	52
4.3	Formation time of persistent teams.	52
4.4	Productivity of persistent team cores.	53
4.5	Dissolution of persistent team cores.	54
4.6	Age composition of persistent scientific cores.	55
4.7	Co-localization of persistent scientific cores.	56
4.8	Interdisciplinary diversity of persistent scientific cores.	57
4.9	Formation time as a function of core knowledge composition.	58
4.10	Career length as a function of core compositions.	58
4.11	Core exclusivity.	59
4.12	Contribution of transient members of persistent teams.	60
4.13	Example of a career of a scientific team.	60
4.14	Temporal dynamics of team success.	61
4.15	Team impact as a function of career stage.	62
4.16	Hot-streaks in team careers.	63
4.17	Team age composition and impact.	63
4.18	Team impact for mono-university and multi-university cores	
	over the decades.	64
4.19	Team productivity as a function of knowledge composition of	
	members	65
4.20	Productivity as a function of core's geographical and age com-	
	position.	66

4.21 4.22	Success as a function of interdisciplinary diversity of the team. Effect of non-core members on team success.	67 68
5.1	Temporal higher-order networks.	74
5.2	Contagion on static and temporal simplicial complexes.	76
5.3	Size of the infected population as function of λ_{\parallel} obtained ana-	
	lytically in the mean-field limit for different values of λ_{\triangle} .	77
5.4	Effect of initial infection size on the onset of the endemic state.	78
5.5	A schematic of temporal simplicial complexes with low and	
	high temporal correlations.	80
5.6	Size of the infected population at the steady state as a function	
	of β_{\triangle} for two different temporal correlations, $\sigma = 0.7$ and $\sigma = 0.3$.	81
5.7	Effect of temporal correlations on the critical group infection	
	rate in higher-order networks.	82
5.8	Effect of heterogeneity in higher-order networks.	83

v

CHAPTER 1

INTRODUCTION

Success is a driving force behind many human endeavors, fueling innovation, progress, and societal advancement at large. Throughout history, achievements born of success in science, arts, and technology have propelled civilizations forward. From Marie Curie's groundbreaking discoveries in physics and chemistry to Leonardo da Vinci's iconic masterpieces, success manifests in diverse forms across disciplines. Something as fundamental as hiring-whether in science, business, or even sports-often relies on clear, quantifiable measures of success. In academia and business, metrics like publications or sales achievements are key, while in sports, evaluating a player's performance is essential to make informed hiring decisions. Besides hiring, governments and NGOs often rely on quantifiable success measures to evaluate public policies and social programs. For instance, metrics like poverty reduction, employment rates, or school enrollment rates are essential in assessing whether programs meet their objectives and justify continued investment. Similarly, in the entertainment industry, success metrics like box office revenue, streaming numbers, and social media engagement shape content creation, casting, and marketing to align with audience interests and profitability. Overall while the definition of success is domain-specific, its quantification is universally important.

Nevertheless, properly quantifying success is a challenging task. In science, it was historically limited due to fragmented bibliometric sources and manual compilations, as noted in the Royal Society Catalogue of Scientific Papers (1867-1925). Only after the digital revolution and the advent of comprehensive databases like Google Scholar in 2004 did unified data sources become widely available. This trend paralleled developments in sports, where datadriven analysis was similarly constrained. The Oakland Athletics' use of sports analytics in the early 2000s, as chronicled in Michael Lewis' book "Moneyball," [1] brought data-driven decision-making in sports to mainstream attention. Secondly, relevant metrics may not be clearly defined or agreed upon, making it challenging to determine what data should be collected. In the arts for example, success is influenced by historical context, market demand, and cultural appreciation [2]. In other words, success may be also associated with recognition, popularity and other collective measures that quantify a community's response to performance [3, 4]. While measures like popularity often correlates with perceived quality, overall quantifying success still remains challenging due to the subjective nature of art and lack of data arising from the opacity of transaction records [5, 6, 7]. Success is not simply the objective performance [8] which represents the measurable achievements and their relationship remains unclear. The distinction between success and performance is evident in terms like 'celebutante,' or 'faminess,' which describe individuals who lack concrete achievements yet are recognized primarily for being famous [9, 10]. Such observations further distance the two concepts [11, 12], increasing the difficulty of pinning down a definition of success.

Recent availability of large-scale data has facilitated advancements in quantification of success. In Science, for instance, the accessibility of multiple, massive data sources such as Web of Science (WoS), Scopus, and OpenAlex [13, 14], along with the collaborative and interdisciplinary effort of scientometricians, natural, computational, and social scientists, has led to the emergence of the interdisciplinary field of Science of science [15, 16, 17, 18, 19, 20], which aims at quantifying the fundamental mechanisms underlying patterns and behaviors in scientific research.

In particular, over the last decade, a large body of literature has focused on the unfolding of individual scientists' careers. The temporal patterns of productivity and impact of researchers were investigated in detail, revealing that the most-cited work of a scientist occurs randomly within her career [21], and that her high-impact papers are likely to come in close succession [22, 23]. Further research explored how the career of scientists is affected by their individual characteristics, such as gender [24, 25] and ethnicity [26], as well as their academic choices and opportunities, including being affiliated with certain institutions [27, 28, 29], moving to a different one [30], switching research topic [31], or dropping out of academia [32]. Funding acquisition is another pivotal for individual scientific careers and their success, supporting crucial aspects such as hiring staff and covering costs of research. However, various factors influence funding success, including disparities faced by women, minority researchers, and those from smaller institutions [33, 34, 27]. Biases against interdisciplinary research exist, while past grant winners tend to receive future funding, perpetuating imbalances in allocation [35, 36]. Elite academic institutions often receive disproportionate funding, fostering exclusive research circles and impacting faculty hiring dynamics [37, 38, 39]. In chapter 2, I provide a new perspective on success in scientific careers, revealing the role of the social network of scientist in funding acquisition. I analyze the careers of ERC winners in Europe and compares them with NSF winners in the US, providing a data-driven comparison of how individual funding success depends on the social network of the scientist. I identify potential hidden biases in ERC grant awards based on the cross-continental collaboration patterns of grant applicants. Moreover, the analyses reveals concerning indicators that research funded by the European Union relies disproportionately on US collaborations. Such a bias risks compromising the independence of European researchers.

Beyond science, temporal career trajectories are also important in Sports. Sports analytics are now commonplace in most major sports, providing clues for individual and team performance to boost success rates [40, 41]. Interestingly, while sports have benefited from scientific methods [1], , they have in turn become a frontier to develop new scientific tools [42] to investigate success, innovation, and learning, as one of the primary domains where growth and success are measurable in a data-driven fashion. Chess serves as a good laboratory to study success as it is a highly intellectual activity that shares similarities to science. Thus, it is often located amid the two domains, a game where players use simple rules resulting in highly complex plays, often developing different personal styles able to influence long-term success in the game. Besides, the volume of online chess games freely available for analyses (several billions), makes chess a perfect candidate for testing hypotheses involving human performance in competitive settings. In chapter 3, I focus on individual careers in the competitive sport of chess. So far chess has predominantly been looked at at the level of single games. Indeed, little attention has been devoted to individual careers and their evolution. I focus on the following question what separates skilled players from the rest? Earlier studies found that the answer is not intelligence [43], and the role of deliberate practice remains heavily debated [44, 45, 46]. quantification of human behaviors and patterns of success. I quantify patterns of success and the effect of skill and career stage on the game openings.

Advancing the discussion from success of single entities, next I focus on teams. Teams are important as many significant advancements are not the result of individual effort alone but rather the outcome of collaborative work of teams. Landmark scientific discoveries such as the identification of the structure of DNA by James Watson and Francis Crick, alongside contributions from Rosalind Franklin and Maurice Wilkins, highlight the power of teamwork in driving scientific innovation. Similarly, the monumental achievement of landing humans on the moon was made possible by the collective expertise and coordination of NASA's Apollo program team. This highlights the importance of going beyond individual careers. Indeed, in the last decades teams increasingly dominate solo authors and in knowledge creation, not only in sheer volume but also in the attention they receive as citations [47, 48]. The increasing predominance and success of teams have been linked to the increasing need for specialized knowledge and interdisciplinary collaborations in solving modernday scientific problems [49, 50]. Advent of technology and the internet plays a major role by allowing communication and thus, collaborations that transverse large geographical distances[51]. Inevitably, composition of scientific teams have garnered increasing interest [47, 50, 52, 53, 54, 55]. More fine-tuned investigation going beyond this simple dichotomy into solo-authors and teams, reveals that the science produced by large teams differs in character from that of small teams. Small teams disrupt science by creating articles which are not often co-cited with its references, while, large teams develop science[56]. Identification of drivers of success in teams becomes crucial. There is not a single driver of success and scientific impact in teams but the picture is complex and multifaceted. Known determiners of impact include team diversity via multiuniversity collaborations [50], inter-member familiarity [57] and prior shared successes of teams [54]. Team diversity in terms of ethnicity, discipline, gender, affiliation, and academic age also has been shown to positively affect impact on the level of disciplines, with ethnic diversity in particular showing the strongest effect on impact [58]. Mechanisms which govern team assembly [47] as well as analytical models which attempts to explain the evolution of team sizes over time on an aggregated level [59] have also garnered some attention.

Despite the vast literature on team science, the temporal patterns of scientific collaborations has so far been largely overlooked and team careers remain an uncharted territory. Akin to those of individual scientists, some scientific teams follow consistent career trajectories, the features of which remain largely unknown. Persistent collaborations among "science buddies" [60] are indeed considered pivotal to scientific research [61], positively contributing to productivity and impact [62, 63]. The few works accounting for the temporal characteristics of collaborations are limited to pairs of scientists [64], or focus on their impact on individuals researchers [62]. In chapter 4, I fill this knowledge gap on career of scientific teams and explore the temporal patterns behind their formation, activity and eventual dissolution. I characterize the composition of persistent cores along multiple dimensions, including age composition, geographical diversity and disciplinary expertise, as well as the role of transient team members. I conclude this chapter by investigating the temporal patterns of success across the career of teams and identifying the compositional correlates of their success.

The exploration of team success in science, particularly through the lens of persistent core and transient team members, leads me to consider how scientific ideas propagate within these temporally dynamic populations. This line of my research connects to broader themes of how social contagions spread through networks, a topic that has received significant attention [65]. More closely related to my interest in understanding the effect of time-evolving groups interaction on contagion, is the sub-field of complex contagions [66] Complex contagion, as defined by Centola and Macy [67], requires multiple sources of activation for transmission. In simpler terms, a single active contact isn't enough to spark adoption; multiple exposures to the same stimulus are necessary for contagion to occur. This concept has been supported by empirical evidence across various contexts and experiments [68, 69, 70, 71, 72]. However, pairwise interactions alone are insufficient to accurately describe social contagion processes such as opinion formation or the adoption of new ideas, which involve intricate mechanisms of influence and reinforcement. A higher-order model of social contagion was proposed [73], where the social system was represented as a simplicial complex, where contagion occurs through interactions within groups of varying sizes. Building on this literature and incorporating the dynamic nature of interactions, in Chapter 5, I model the spread of ideas across populations with evolving structures using an agent-based approach to social contagion. Recognizing that human interactions vary over time, I introduce temporally evolving population structures represented by temporal simplicialcomplexes. Using simulations, I reveal the crucial role of persistent groups and temporal correlations in the contact structure in the spreading of innovations. By studying how ideas flow through groups that vary in their stability and member composition, we can uncover strategies to accelerate the diffusion process and potentially enhance the rate of collective scientific discovery.

Overall my investigations of successful careers of individuals and teams in a variety of domains add to the scientific conversation on the emerging field of Science of Success, opening directions for further investigations on identifying the paths to success in scientific and sports careers.

CHAPTER 1. INTRODUCTION

CHAPTER 2

ROLE OF COLLABORATION NETWORK IN SECURING INDIVIDUAL SCIENTIFIC FUNDING

2.1 Introduction

Thousands of scientists regularly submit grant applications to secure funding and develop their most innovative ideas and propel their careers forward. In a scientist's career, funding is required for multiple aspects, including hiring postdocs and PhD students, renovating labs by acquiring new instruments and equipment, and covering publishing fees. Hence, it is crucial to understand what factors influence funding success of a scientist. However, selection biases have been documented to impact different groups of researchers. Women [33], researchers from minority groups [34] and from small institutions [27] are known to have lower grant success rates. A funding bias has also been observed against interdisciplinary research [35], possibly due to the difficulty of satisfying the requisites of multiple scientific communities. Moreover, past grant winners are more likely to win future funding [36], suggesting the existence of a Matthew effect that potentially reinforces small initial imbalances in funding allocation leading to large disparities later in the careers. At the university level, closed circles of elite academic institutions often receive disproportionate amounts of funding [37], giving rise to rich clubs [38] of researchers who collaborate predominantly among themselves, eventually affecting the dynamics of faculty hiring processes [39].

The ERC and the NSF are among the major research funding agencies in Europe and the US, providing crucial support in the establishment and consolida-

tion of successful careers across a variety of scientific domains. Within the EU, the funding landscape is heterogeneous [74], with a preference towards countries that can retain their scientists and attract outside talent [75]. However, little is known about the intertwined nature of the European and American funding landscape.

In this chapter, I focus on careers of ERC and NSF awardees, and determine the role of their social network on funding acquisition. I detect asymmetries in the distribution of cross-continental collaborations among top-funded researchers in Europe and the US, discovering a consistent pattern of collaboration in the careers of successful ERC grant winning scientists.

2.2 Data

The data from openalex.org used in this chapter is openly accessible for download using the API https://api.openalex.org/works. Details of awardees were downloaded for ERC from https://erc.europa.eu/project-statistics/projectdatabase and for NSF from https://www.nsf.gov/awardsearch/download.jsp. I considered the careers of ERC winners since 2008, and compared them to those of NSF awardees over the same timeframe. As the countries eligible to receive ERC awards have changed since the launch of the funding program in 2008, I narrowed the analysis to 6,260 ERC winners from 2023-eligible countries [76]. In particular, I used disambiguated authors from the OpenAlex dataset [13] an open and comprehensive catalog of scholarly papers, authors, institutions and related data [77].

2.3 Cross-continental collaborations of ERC and NSF awardees

As a first step, I quantified the relevance of US/EU collaborators for ERC/NSF awardees as a function of their academic age, measured from the date of their first publication (Fig.2.1).

Specifically, by measuring the fraction of cross-continental collaborations per paper, I found that ERC winners have the largest share of US-based collaborations early in their career, with a sharp decline after about 7 years. By contrast, no significant temporal trend is observed for NSF awardees, whose fraction of European collaborators remains approximately constant over the first twenty years of their careers. I obtained consistent results after splitting ERC winners by grant type, finding that winners of ERC Starting Grants — awarded on average 11.3 years after their first publication — have the highest share of US col-



Figure 2.1. Asymmetric collaborations of European and American funded research.

Mean percentage (error bars represent standard error) of cross-continental coauthors per paper by ERC awardees (green line) and by NSF winners (purple line) and as a function of career age.

laborations before winning the grant (20.9%), followed by Consolidator (13.6%, 14.6 years) and Advanced (13.3% for, 24.8 years) awardees.



Figure 2.2. Collaboration imbalances between EU based scientists and US based scientists.

EU-based scientists, who exclusively operate within the EU, exhibit considerably fewer collaborations with their US counterparts compared to ERC winners.

In order to assess whether this imbalance is simply a byproduct of a different collaboration culture between EU and US, I checked the trends of cross continental collaborations for all EU and US based scholars. To test this hypothesis, I have considered two groups of scientists. Operationally, two strict definitions have been considered, where I defined as EU based scientists all researchers in the OpenAlex dataset with at least 20 papers and having only being affiliated to the EU, and as US based scientists all researchers with at least 20 papers and having only being affiliated to the US. I evaluated the percentage of crosscontinental collaborations among EU-based scientists by compiling the papers of all scientists in the considered cohort published in the *n*th year of their career. I then calculated the average percentage of US coauthors per paper across the collected papers. This procedure is performed for various career ages, divided into 1-year age bins. The same procedure is applied to the US-based pool to compute the corresponding measure. Looking at the percentage of crosscontinental coauthors for EU and US based researchers as a function of time, I find that scientists who are and have been exclusively based in the EU collaborate notably less with their US counterparts compared to ERC winners (Fig.2.2). Moreover, US based scientists tend to have more cross-continental collaborators than EU based ones, a trend which is stronger later on in their career. These findings — opposite of what observed for ERC and NSF awardees — suggest that the patterns of cross-continental collaborations observed for grant awardees are not a simple consequence of general unbalanced patters of collaborations between the EU and the US, and that the award of ERC grants is linked in a nontrivial way to the existence of scientific links with the US research ecosystem.



Figure 2.3. Differences in collaboration patterns before and after winning the grant.

Mean percentage (error bars represent standard error) of cross-continental coauthors per paper by ERC awardees (green line) and by NSF winners (purple line) and as a function of career age.

2.4 US collaborations before and during an ERC grant

Investigating the asymmetry between ERC and NSF winners further, I compared the fraction of US collaborations before and after winning the grant in a five-year time window (Fig.2.3). The results showed that, before winning the ERC grant, EU scientists tend to include significantly more US based coauthors (194.0% more than EU collaborations by NSF). However, these cross-continental collaborations typically wane after receiving the award, as the US collaborator fraction decreases during the grant period (101.0%). On average, there are 208.8% more US collaborations before winning the grant than during the grant period. By contrast, I found no significant change in number of EU collaborations for NSF awardees. These findings suggest the presence of potential selection biases when awarding ERC grants to early career EU scientists that have a strong network of collaborations in the United States. After winning the grant and securing a stronger academic stability and independence, the weight of US based collaborations typically starts to decay for EU researchers.

Given the ERC's aim to foster academic independence, an alternative possible explanation for the decline in US collaborations among ERC winners postgrant a trend could be an overall post-award collaboration decrease. To investigate such a possibility, I computed the number of ERC winner coauthors before and during grants (in both cases I considered 5 year windows). Then, I measured the percentage change in the number of unique coauthors during the grant compared to before (Fig.2.4).



Figure 2.4. *Number of collaborators before and during an ERC grant. Distribution of percentage change in number of unique coauthors during ERC grant as compared to before winning the grant. Vertical black line marks the positive mean of the distribution. Most ERC winners expand the set of their collaborators after receiving the award.*

I found that majority of ERC winners explore more diverse collaborations after winning, instead of exploiting fewer collaborators. This suggests that the decrease in post-grant US collaborations is not associated with a reduced openness of ERC awardees to collaborate with colleagues.

Each year, a significantly greater number of scientists receive funding from the National Science Foundation as compared to the European Research Council. Thus, it is important to verify whether the imbalances in cross-continental collaborations before and after the grant award are due to the larger bases of NSF funded researchers. To test such an hypothesis, here I limit the analysis to the top NSF recipients, selected to match the number of ERC winners. Specifically, for each year from 2008 onwards I select the top 'x' NSF awardees based on their ranking by citation counts or alternatively by the grant award amounts, where 'x' is the number of recipients of ERC award. My findings remain qualitatively unchanged with both criteria, indicating that the collaboration imbalance between the European Union and the United States persists (Fig.2.5).

While imbalances between ERC and NSF awardees persist, it is worth noting that elite NSF winners do exhibit a greater tendency to engage in international collaborations with EU scientists. This interesting outcome merits further investigation in future research.



Figure 2.5. Comparison between ERC winners and top NSF awardees.

When compared to those authored by top-cited NSF researchers, EU funded papers had 99.7% more cross-continental collaborators before grant award and 35.3% after the award. When selecting top NSF researcher in terms of grant award amount, EU funded papers had 142.0% more cross-continental collaborators before grant award and 78.1% after the award.

2.5 Pre-award mobility of ERC vs NSF awardees



Figure 2.6. *Asymmetric mobility of European and American funded researchers. Box plots of the number of distinct universities appearing in the papers' affiliations of ERC (green boxes) and NSF (purple boxes) winners conditioned on the academic age at the time of grant award. Horizontal black lines represent the median of the distributions.*

Furthermore, my results indicated a higher mobility of ERC awardees before winning the award when compared to their NSF counterparts. I found that early career stages for both EU-based and US-based awardees are characterized by an increasing mobility, but visible differences appear after the first five years of their career (a typical length of a doctoral degree in many countries). Specifically, the median number of distinct universities appearing in the affiliations reported on ERC awardeess papers is consistently higher than affiliations per NSF awardees of comparable academic age (Fig.2.6). Considering the number of different countries, instead of different affiliations, yielded similar results (Fig.2.7).

The greater number of countries visited by ERC awardees might in principle be a byproduct of the different countries within the EU. To test for this hypothesis, I have produced a similar analysis where I regard the European Union as a unified entity and I do not consider traveling between its constituent countries as a form of mobility. The main observation remains unaltered, with recipients of European Research Council grants exhibiting greater levels of mobility than recipients of National Science Foundation grants (Fig.2.8).

Thus, mobility patterns between the EU and the US exhibit an imbalance similar to the one I observed for collaborations. This suggests that the markedly different career trajectories and degrees of international experience of ERC and



Figure 2.7. *Mobility of scientists across countries before grant award. Comparison of mobility measured as number of distinct countries appearing in papers' affiliation*

before winning the grant. Dashed horizontal black lines indicate the mean values.

NSF awardees may have implications for their research perspectives. Indeed, although a substantial proportion (49.7%) of ERC winners had prior affiliations with US institutions, NSF winners were considerably less likely (only 13.5%) to have had previous affiliations in any EU-based academic institution.



Figure 2.8. *Mobility differences before grant award: EU as a single state.* Comparison of mobility before winning the grant between ERC and NSF recipients if EU is regarded as a unified entity. Dashed horizontal black lines indicate the mean values.

2.6 Wider collaboration patterns of ERC and NSF funded researchers



Figure 2.9. Ratio between percentages of ERC and NSF collaborations with EU countries.

A value greater than 1 indicates that ERC winners exhibit a higher preference for collaborating with researchers from the country in question, while a value of less than 1 value suggests that NSF winners show a greater inclination towards collaborating with scientists from this country.

Subsequently, I aimed to quantify the likelihood of finding collaborators across different countries for ERC and NSF winners, as a measure of embeddedness within the ERC and NSF funding landscape. Thus, I computed the ratio between the percentage of ERC and NSF winners' collaborators that belong to a given EU country(Fig.2.9). As expected, EU countries contribute more to the collaboration network of ERC awardees than NSF ones. On average, EU countries produce 8.4 times more collaborators to ERC awardees than to NSF ones. Large countries such as Germany, France, Spain, and Italy are predictably among the top ERC-embedded EU nations. Central european countries such as Belgium and the Netherlands, also hosting EU institutions, emerge as leaders, with the highest overabundance of collaborators for ERC awardees.

I then broadened the analysis to non-EU countries, finding that Taiwan, China, Georgia, South Korea, and Singapore, host a greater fraction of collaborators for NSF winners compared to ERC ones. Countries such as Switzerland, United Kingdom, Serbia and Russia attract instead more ERC collaborations (Fig.2.10). The case of the UK is particularly interesting, as UK-based scientists had received approximately 20% of all ERC awards prior to Brexit. The find-



Figure 2.10. Ratio between percentages of ERC and NSF collaborations with non-EU countries.

A value grater than 1 indicates that ERC winners exhibit a higher preference for collaborating with researchers from the country in question, while a value less than 1 suggests that NSF winners show a greater inclination towards collaborating with scientists from this country.

ings revealed that the UK exhibits only weak ties with EU-based ERC awardees, while maintaining strong connections with NSF researchers based in the US. Indeed, the UK only produce 1.8 times more collaborators to ERC awardees than NSF ones, against an EU average of 8.4 times, outranked only by Bulgaria. A similar situation is observed for Switzerland.

2.7 Universality of the EU-US imbalance

Finally, I explored the cross-continental EU-US imbalance in collaboration patterns across different scientific disciplines as per OpenAlex Dataset (Fig.2.11). The analysis suggested that the imbalance is present across the top ten largest fields for both ERC and NSF winners. Geology (244.6% more US collaboration by ERC winners), Medicine (208.5% more) and Biology (189.7% more) are the most unbalanced fields, while Physics (99.4%), Mathematics (73.3%) and Materials Science (46.0%) the least unbalanced ones. I also measured the fraction of cross-continental collaborations across the different ERC-specific domains, and how they changed before the award and during the grant period (Fig.2.12). I grouped winners based on the issuing ERC panel [78] into Physical Sciences and Engineering (PE), Life Sciences (LS) and Social Sciences and Humanities (SH). Overall, LS winners have the highest percentage of US collaborators (20.1% be-



Figure 2.11. *Discipline-wise quantification of asymmetrical cross-continental collaborations patterns for ERC and NSF awardees.*

fore and 11.6% during the grant), but also the largest decrease once funding have been awarded (-42.2%). PE (15.2% before and 10.2% during) awardees tend to rely less on US collaborations than SH (15.8% before and 12.2% during), particularly so after the grant award (-32.9% and -26.4% respectively).

At the level of the individual panels, "Individuals, Markets and Organisations" (SH1) records the highest fraction of US collaborations (30.3%), while "Synthetic Chemistry and Materials" (PE5, 11.2%) the least before securing funding. Awardees from almost all panels display a drop in US ties after winning the grant, with "Physical and Analytical Chemical Sciences" (PE4) displaying the largest decrease (-50.0%). The only exception appears to be "The Study of the Human Past" (SH6, 9.4% increase). In physics, "Universe Sciences" (PE9) awardees have the tightest links with the US (20.3% of collaborators) before the grant is awarded, and display the small decrease in cross-continential collaborators afterwards (-18.0%). "Fundamental Constituents of Matter" (PE2) and "Condensed Matter Physics" (PE3) winner tend to have fewer US collaborators before award (16.4% and 14.5%), and among the largest drops in such numbers among all panels as a consequence of receiving an ERC (-33.5% and -42.2%).

2.8 The case of physics

Collaboration and mobility imbalances

As a case study, I analyze collaboration and mobility imbalances for researchers recognized as physicists as per the OpenAlex dataset. In particular, I select all



Figure 2.12. Percentage of US coauthors per paper for ERC winners from different grant panels before and after winning the grant.

researchers who score above 90 (over 100) for physics in the OpenAlex dataset, leading to 1295 ERC awardees and 10249 NSF awardees for physics (19.0% of the whole populations. Fig.3a on imbalances in cross-continental collaborations was produced with the same choice of the threshold for disciplinary categorization. I note that researchers may be assigned to multiple disciplines).

The dataset includes 85.0% of the awardees from PE2, 79.2% from PE3, and 87.0% PE9, the three main ERC panels for physics and astrophysics. The dataset also includes 607 ERC winners recognized as physicists awarded from different panels, recognizing the interdisciplinary nature of some physics topics, applications of physics to other domains, and the interdisciplinary nature of some ERC panels (e.g. PE8, PE11). Results reported below are not significantly affected by the exact choice of the threshold for physics categorization in the Open Alex dataset.

I observe similar trends compared to the full population of grant recipients. However, cross-continental collaboration imbalance for physicists is generally less pronounced. The maximum imbalance is still observed at the initial stages of an academic career, in particular around the sixth year (Fig.2.13a). Physicists who received funding from the European Union (EU) exhibited a 112.7% increase in cross-continental collaborators prior to receiving the grant, compared to a higher value of 194.0% for all ERC winners. After receiving the grant, this percentage decreased to 51.1% for physicists, in contrast to 101.0% for all ERC winners (Fig.2.13b). As for the general population of researchers, ERC awarded physicists are more mobile than their NSF counterparts (Fig.2.13c).



Figure 2.13. Co-authorship and mobility differences between funded physicists in the EU and the US.

a) Mean percentage (error bars represent standard error) of cross-continental coauthors per paper by ERC (green lines) and NSF (purple lines) physics awardees (solid lines) as a function of career age, compared to the whole population of grant recipients (dashed lines). Physicists exhibit a comparatively reduced collaboration imbalance in comparison to the overall group of winners. *b)* Average percentage of cross-continental coauthors per paper by ERC winners' (green bars) and NSF awardees (purple bars) in physics before grant award and during the grant period. Vertical black lines represent standard errors of the values. *c)* Box plots of the number of distinct universities appearing in the papers' affiliations of ERC (green boxes) and NSF (purple boxes) winners for physics at the time of grant award.

Country-wise collaboration patterns for physics awardees

In the previous section, I quantified the fraction of collaborators across different countries for ERC and NSF winners, as a measure of embeddedness within the ERC and NSF funding landscape. I find that similar trends is observed if considering only physicists, both within (Fig.2.14a) and outside (Fig.2.14b) the EU. Countries like Cyprus or Lithuania appear to be more embedded in the ERC collaboration network when focusing on physics only.

To summarize, focusing on physics, similar imbalance trends persist, with a similar peak of US collaborations for ERC winners around the sixth year of their careers, although less pronounced (Fig.6a). Physicists funded by ERC had 112.7% more US coauthors than vice versa pre-grant, while post-grant this percentage tapered to 51.1% (Fig.6b). ERC winners in the physics domain are also more mobile than their US counterparts (Fig.6c).

2.9 Discussion

In this chapter, we tracked careers of individual funded scientists to reveal the role of scientists' social network in successful career. Our results shed new light on the ties between the EU and the US academic ecosystems. The findings highlight the unbalanced relationship between the two, with early collaborations



Figure 2.14. Countrywise collaboration patterns of ERC and NSF winning physicists.

a) Ratio between percentages of ERC and NSF collaborations with EU countries. A positive value indicates that ERC winning physicists exhibit a higher preference for collaborating with researchers from the country in question, while a negative value suggests that NSF winning physicists show a greater inclination towards collaborating with scientists from this country. *b)* Ratio between percentages of ERC and NSF collaborations with non-EU countries.

with US institutions providing a crucial advantage for success in securing top EU funding.

Such imbalance might compromise the independence of European researchers, in particular during early career stages, by pressuring them to align with the interests of the US scientific community in order to secure academic success [79]. Moreover, such selection bias in EU funding might be associated with a number of possibly unnecessary collaborations or career moves, motivated more on strategic rather than scientific grounds. In fact, according to my analysis, researchers willing to receive major funding from the EU are expected to be more mobile and go through a wider diversity of relocating experiences than their American counterparts. In particular, landing a job in the US appears as a key factor to secure an ERC award later on, while I found no evidence of the opposite.

Overall, the results draw a worrying picture of the potentially subordinate role of the European research community with respect to the American one. Such a claim would already be supported by independent studies on the hiring market [80], which showed that American top universities hire a substantial fraction of researchers from European universities that score far below than them in international rankings. It is possible that the influence of US collaborations might be mitigated in the future years, following the recent changes in ERC evaluations guidelines, which give more weight to the scientific content of project proposals over applicants' resumes [81].

The analysis does not explicitly account for NSF's exclusion of certain disciplines, such as medical sciences, which are typically sponsored by NIH. Addi-

20

tionally, differences in grant durations between ERC and NSF awards, ranging from a few months to five years for NSF and typically five years for ERC, are not factored into the results. While our analysis of the top 10 fields demonstrated that the imbalance in US-EU collaborations is likely universal, it is possible that this imbalance is absent in some of the fields funded by NIH, a potential result which would merit further investigation. Furthermore, the variation in grant durations may influence collaboration dynamics, with longer ERC grants fostering sustained, stable partnerships which might be harder to maintain across continents explaining the decrease in US collaborations after winning an ERC grant. To address these issues, future research could employ a matched pair analysis, pairing researchers based on grant duration. This approach would help control for structural differences between the two funding systems. Expanding the analysis to include NIH-funded projects could facilitate a more fair comparison of the US ecosystem with the EU.

Code availability

The code used in this chapter is available at https://github.com/chowdhary-sandeep/NSF_vs_ERC.
CHAPTER 3

QUANTIFYING PERFORMANCE IN CHESS CAREERS

3.1 Introduction

In the previous chapter I have discussed quantitatively how social networks help in developing a successful careers in science. Beyond science, underpinning quantitatively the drivers of success is a key problem in sports. Indeed, the recent availability of large-scale datasets is nowadays providing an unprecedented opportunity to study the drivers of human performance in all such different domains. Analyzing individual performance with data has become essential in sports. The impact of this approach was evident in baseball with the transformative story depicted in "Moneyball" [1]. The realization that statistical data can inform strategies and player evaluations revolutionized the game. Similarly, in tennis, network techniques revealed Jimmy Connors as a standout player from the past [42]. Today, sports analytics are widespread, offering insights to enhance both individual and team performance [40, 41]. Notably, sports have not only benefited from scientific methods but have also become a testing ground for developing new scientific tools to study success, innovation, and learning in a data-driven manner.

In this chapter, I focus on individual careers in the competitive sport of chess, extending my exploration of success in individual careers in this thesis. Chess, as a domain bridging intellectual pursuit and competitive sport, provides a unique model for studying performance and career dynamics. Thus, it is often located amid the two domains, a game where players use simple rules resulting in highly complex plays, often developing different personal styles able to influence long-term success in the game. Besides, the volume of online chess games

freely available for analyses (several billions), makes chess a perfect candidate for testing hypothesis involving human performance in competitive settings. So far chess has predominantly been looked at at the level of single games. For example, past research focused on the role of memory in games [82] and showed that opening popularity follows the well-known Zipf's law [83]. However, these analyses did not use individual player-level data, treating games from different players on equal footing[82, 83], or focused on a small number of players [84]. Indeed, little attention has been devoted to individual careers and their evolution. In particular I ask—*what separates skilled players from the rest?* Earlier studies found that the answer is not intelligence [43], and the role of deliberate practice remains heavily debated [44, 45, 46].

Here I perform a comprehensive large-scale analysis of the habits of skilled and less skilled individual players over time, providing an anatomy of human performance in the popular game of chess. I characterize players' careers in terms of hot-streaks, diversity and specialization in the opening sequences of their games, and analyze their diversity as a function of career stage. I find evidence for the presence of both hot and cold streak phenomena, revealing a surprising tendency for beginners to have longer hot-streaks as compared to expert players. By sequencing the opening moves of players at different skill levels, I show that beginners start with more diverse set of first moves, while advanced players and experts rarely start their games differently when playing as white. Yet, expert players display a broader response repertoire, showing the ability to surprise their opponent with a greater variety of responses. Moreover, when accounting for different variations of the openings, experts show a deeper knowledge of different variations within the same line, hinting at a deeper understanding of the game. Lastly, analyzing behaviour in time, I find that players explore more during the beginning of their careers, but tend to specialize using and exploiting only fewer openings at later career stages. Overall, this large-scale characterization of individual gaming behavior supports chess as a suitable laboratory to quantitatively investigate individual careers and human performance, demonstrating simple differences in playing habits and behaviours of beginners and experts.

3.2 Methods

Data

Lichess provides an extensive game database, and my analysis is based on its publicly available dataset. The data is openly accessible for download from https://database.lichess.org/. I use all games played on the online chess server from lichess.org between 2013 and 2016. There are different games available to the players on the platform: bullet, blitz, and rapid. The analysis presented in this chapter is restricted to Blitz games, which are fast and tactical but still allow for some strategy in the game overall unlike bullet games which last only 1 minute at most and littered with pre-moves. The most popular time controls for blitz are 5 mins and 3 mins. I specifically focused on this type of games since "speed" chess is played across all levels, from beginners to grandmasters. The most popular blitz time controls are 5 minutes and 3 minutes. Regarding data cleaning, I applied a filtering criterion where players with fewer than 100 games in their career were removed. This choice was motivated by the observation that player ratings fluctuate significantly during their first 100 games before stabilizing. While this early career phase may be interesting in its own right, it introduces additional variability that is beyond the scope of this study and may be the target of future investigations. Our final dataset includes 123 million games played by 0.98 million players.

Matching on lichess.com

Overall, Lichess.com offers three primary pairing methods to cater to different player preferences: Quick Pairing, where players select preset time controls like 3+0 or 5+3 and are automatically matched with opponents of similar ratings, typically within a ± 100 to ± 150 point range; Arena Tournaments, which prioritize fast matchmaking by pairing players based on availability and similar rankings, allowing multiple games within a set duration but occasionally leading to repeat opponents; and Swiss Tournaments, a structured format where players are paired based on their current scores and tie-breaks, ensuring competitors face others with similar performance levels while avoiding repeat matchups, with all players waiting for each round to conclude before proceeding to the next, promoting fairness and a standardized competitive experience. Potentially there is a way to play friends repeatedly on lichess.com who can act as teachers– a motif which can be explored in further research with this data. I will add these insights in the methods section, as well as limitations in the discussion section.

Measuring opening diversity

I measure the diversity of openings of a player by calculating the Shannon entropy [85] of the distribution of frequency of opening moves or opening sequences (see Fig. 3.12, 3.13 and Fig. 3.16 respectively). Note that for the analysis in Fig. 3.13, I selected only games where the player starts as white.

Null models for hot and cold streaks

To calculate the expected lengths of hot streaks in a player's career, I build a null model where I reshuffle the temporal order of the player's games, but preserving the total number of victories, losses and draws. Such shuffling of the order of games within each player's career breaks the temporal correlations between game outcomes[22].

Then, I compute the length of each hot and cold streak (sets of consecutive wins and loses) observed in this reshuffled sequence. The presence of hot- and cold-streak phenomena can be then investigated by comparing the number of hot and cold streaks of a given length ℓ in the actual careers with respect to these reshuffled sequences.

Chess concepts

Openings

A chess opening is the initial stage of a chess game—a sequence of first few moves. It usually consists of established theory; the other phases are the middlegame and the endgame. All games can be associated with a unique main opening line, within which there can be many variations. Many opening sequences have standard names such as the Sicilian Defense, Ruy Lopez, Italian Game, Scotch Game etc.

Chess rating systems

I present a very short overview of different rating approaches developed for chess. This analyses are based on the Glicko-2 rating systems.

Elo system, invented by Arpad Elo, is the most common rating system for chess. It is used by FIDE, other organizations and some Chess websites such as Internet Chess Club and chess24.com. Games are scheduled by matching together players of similar ratings.

Glicko-1 system[86] is a more modern approach, invented by Mark Glickman as an improvement of the Elo system, which preserves the philosophy or the Elo rating approach while making it more accurate. In the Glicko system, a player's rating not only changes due to game outcomes, but also from their "ratings deviation", which measures the uncertainty in a rating due to both game outcomes and also from the passage of time when not playing. At the cost of being more mathematically complex, the Glicko rating system is known to have a better prediction accuracy than Elo, and it is a popular choice for new games and sports. The Glicko system has a initial rating starting at 1500. It



is used by Chess.com, Free Internet Chess Server and other online chess servers.

Figure 3.1. Distribution of player ratings in data.

Distribution of Glicko-2 ratings averaged over the career of a player separately for the different time controls i.e. Bullet, Blitz and Rapid.

Glicko-2 system is a refinement of the original Glicko system and is used by Lichess, Australian Chess Federation and other online websites. It achieves even better accuracy by controlling for volatility. Volatility measures the degree of expected fluctuation in a player's rating– it is low when the player performs at a consistent level and high when a player has erratic performances (e.g., when the player has had exceptionally strong results after a period of stability). I show the Glicko-2 ratings of all players in my dataset in Fig. 3.1. Initial ratings start at 1500. Here, I associated to each player its rating averaged in all their games (to account for common early fluctuations, for each player I do not consider their rating in the first 100 games). The career lengths of players as a function of their skill is shown in Fig. 3.2. As shown, on average experts tend to play more games than beginners. I note that players of similar ratings are matched to compete. However, small rating differences among matched players might persist.

Separating players by skill level

I separated the player into the 4 skill levels as follows. I first arranged the players in ascending order of their *Glicko-2* rating (average calculated over all their games). I then created Glicko-2 rating bins that divide players in 4 equally sized skill categories. Finally, I labelled these bins as—*beginner, intermediate, advanced, expert* respectively.

Opening variations

For a given chess opening there are multiple variations, as players can explore different moves after the main opening line is established. For example, the *Sicilian Defence* begins with the following moves 1. e4 c5. The *Sicilian Defence: Najdorf Variation* of Sicilian is 1.e4 c5 2.Nf3 d6 3.d4 cxd4 4.Nxd4 Nf6 5.Nc3 a6, while the *Sicilian Defence: Dragon Variation* is 1.e4 c5 2.Nf3 d6 3.d4 cxd4 4.Nxd4 Nf6 5.Nc3 a6.

3.3 Careers in chess

In this chapter, I rely on large-scale data extracted from *lichess.org*, a popular open-source Internet chess server, consisting of 123 million games between 0.98 million players (see *Methods*). In the *lichess* dataset, each player's career can be tracked over time, with detailed information on each of the played games, i.e. moves, opening, win/loss, and its skill level. This is quantified by the *Glicko-2* rating (see *Methods* for a detailed discussion of different ways to measure skill in chess), which measures the level of past performance of the player, it increases when a player beats an opponent and decreases upon a loss. As an illustrative example, in Fig. 3.3 I show the career of Grandmaster (GM) Magnus Carlsen on *lichess.org*, indicating his *Glicko-2 rating* in each game and the game outcome: win, loss or draw.

3.4 Hot and cold streaks in chess careers

Figure 3.3 suggests that for *GM Carlsen* wins and losses tend to be clustered together. Indeed, prior works tracking wins and losses in sports hotly debate



Figure 3.2. *Distribution of the number of games per player. Distribution of the total number of games by a player for the 4 skill categories.*



29

Figure 3.3. Visualization of the career of the Grandmaster (GM) Magnus Carlsen.

Wins and losses of GM Carlsen drive the rating up or down.

the existence of hot-streaks [87, 88], a phenomenon that has also been found to be ubiquitous in artistic and scientific careers [22, 23].

Long streaks of chess wins are reminiscent of players entering the so-called *zone*, a state of focus where peak performance is possible [89, 90]. **Detecting hot-streaks by comparing against a shuffling model.** To quantitatively check for the existence of such phenomena in all chess careers, I calculate the length of hot (series of wins) and cold (series of losses) streaks for each player in the dataset, and compare them with lengths expected in a null model for each player which shuffles the temporal order of the player's games, thus washing out temporal correlations in game outcomes (see *Methods* for details). In Fig. 3.4 I show the resulting curves, properly normalised with the null model. I find the existence of statistically significant hot streaks, possibly associated with confidence spillovers from previous victories. Long streaks of chess wins are reminiscent of players entering the so-called *zone*, a state of focus where peak performance is possible [89, 90].

Autocorrelation test for hot-streaks. Here the definition of hot-streaks is based on a hard constraint, where one loss ends a hot-streak. A more loose statistical way of quantifying hot-streak phenomena is via measuring autocorrelations, a method which has been used for detecting hot-streaks in literature [91, 92]. I compute auto-correlation with a lag of one game among the outcomes and find it to be positive for all skill levels (Fig. 3.5). This finding suggests that consecutive games have similar outcomes, thus further confirming the presence of hot-streaks via.



Figure 3.4. Hot and cold streaks in chess careers.

Relative number of hot streaks (red) and cold-streak (blue) of length $\ell \ge L_{streak}$ as a function of L_{streak} calculated for each player. Results are averaged over all players. Losses tend to be more clustered than victories as individual cold streaks

Role of rating advantage over opponent in hot-streaks

Observed hotstreaks could also occur simply due to consecutively facing weak opponents or large time-gaps between consecutive games. Players are selected to play together if they have similar rating scores. However, small rating differences might still exist, and impact the length of hot streaks, which could be positively influenced by consecutively facing weaker opponents. To test this, in Fig. 3.6 I plot the average rating advantage over the opponent in hot streaks as a function of their length. Results are averaged by considering all the games composing a given hot streak of a certain length, and over all hot streaks of that length. I find that longer hot streaks are associated to a higher average advantage over opponent, thus supporting the idea that weaker opponents might cause hot-streaks. Thus I find that the length of hot-streak and the average rating-advantage over the opponent in the streak are slightly correlated



Figure 3.5. Autocorrelation test for hot-treaks.

Distribution of autocorrelation in game outcomes for a player's career with lag 1 aggregated into 4 skill categories.



Figure 3.6. *Role of rating advantage over opponent in hot-streaks. Average rating advantage over opponent during the hot-streak as a function of hot-streak length* $l_{hotstreak}$.

($\rho_{spearman} \sim 0.28$), explaining only partially the observed behavior.

In Fig. 3.7, I also compare the rating advantage over the opponent who broke the hot-streak with respect to the one of the opponent in the preceding game (the last one composing the hot-streak) as a function of hot-streak length. The relative advantage is negative, which implies that opponents who break a hotstreak are consistently stronger than opponents who got beaten during the hotstreak. However, I note that this result is somehow expected, as in the chess gaming platform analysed in this chapter, a player who keeps winning consistently (in a hot streak) will likely be matched with higher-rated opponents in the next game.

Role of time-difference between successive games in hot-streaks

I also investigated the effect of time-difference between games and possible breaking of hot-streaks due to time gaps. In Fig. 3.8 I plot the difference between the time-gap before the game where the hot-streak was broken and the time-gap before the preceding game as a function of hot-streak length. For short hot-streaks (lengths 2,3,4) I find a positive relative time-gap before the streak ending game, hinting that such hot streaks could also be disrupted by player taking a break before the next game. The overall trend remains unclear for longer streaks, which are the ones which are really associated with the so-called hot streak phenomenon. I quantify this effect by performing a Spearman's correlation test between the time-gaps before the streak breaking game and the length of the hot-streak. I found a significant but very weak correlation ($\rho_{spearman} \sim 0.05$), hinting that timegaps between games do not play a major role in breaking hot-streaks.



Figure 3.7. Rating disadvantage in the hot-streak breaking game.

Average rating advantage in the streak ending game compared to the advantage in preceding game as a function of hot-streak length $\ell_{hotstreak}$. On average, players have a rating disadvantage in the game where they finally lose after a hot-streak.

Effect of repeated matches with the same opponent on hotstreak behaviour

In other sports, it is well known that players need to vary their strategies against an opponent, so that habits are not exploited for strategic advantage. However, such variation in strategy is not very relevant in online chess games, where repeated matches are rare. Fig. 3.9 shows the frequency of repeated matching of players. I find that 81% of player matches never repeat, and only 1.6% player pairs play more than 5 games with each other. Thus, repeated games do not significantly affect these results.

In addition, to specifically check what happens when opponents are repeated, I investigate the balance of wins and losses in repeated match-ups among two players (with at least 20 games between them). In Fig.3.10, I show the balance of outcomes, computed as $\left|\frac{n_{wins}}{n_{games}} - 0.5\right|$. I find that while a majority of repeated match-ups lean towards balance (50-50), there exist some pairs which are highly imbalanced, with one player dominating and winning consistently, which might result in some hot-streaks. However, as stated earlier, these repeated games represent only a tiny fraction of all the games. Thus my results are not affected by repeated matches against the same opponent.

3.5 Effect of player skill on likelihood of hotstreaks.

I refine the analysis of hot-streaks by further separating players by skill (Glicko-2 rating). Categorizing players into 4 categories - *beginner, intermediate, advanced,*



Figure 3.8. *Role of time-difference between successive games in hot-streaks Relative timegap before streak-ending game– calculated as the difference of timegap before the streak ending game compared to the timegap in preceding game as a function of hot-streak length* $\ell_{hotstreak}$.

expert (see *Methods*). I find that weaker players experience comparatively longer hot streaks than stronger players (Fig. 3.11). A reason for this could be that confidence spillovers from last victory may have greater impact on future outcomes at a lower skill levels.



Figure 3.9. Distribution of number of repeated matches in the dataset.



Figure 3.10. *Distribution of balance in outcomes in player match-ups for the 4 skill categories.*

Distribution of balance in outcomes in player match-ups for the 4 skill categories. A match-up is balanced (=0) if the players had as many wins as losses, and completely imbalanced (=0.5) if one players always beats the other.



Figure 3.11. Hot-streaks and player skill.

Relative number of hot streaks of length $\ell \ge L_{hotstreak}$ as a function of $L_{hotstreak}$, averaged over the players in each skill categories separately (i.e. beginner, intermediate, advanced, expert). Weaker players have longer hot streaks than more expert ones.



Figure 3.12. Diversity and specialization in the first move of the game.

Boxplots showing diversity (entropy) of first move by a player as white, calculated over all players individually and aggregated into the 4 different skill levels. Weak players start games with diverse collection of first move as white when compared to stronger players.

3.6 Specialization in opening play and skill level

Another possible driver of the observed disparity in hot streaks across beginners and experts can reside in how experts diversify their moves. In competitive sports, some players diversify their techniques while others may specialize. Strategy diversification might make players harder to predict, thus enabling them to surprise their opponents. By contrast, specialization, e.g., deeper knowledge of certain opening positions, may allow players to exploit opponents navigating familiar situations. Indeed, such an exploitationexploration (specialization-diversification) dichotomy is a common mechanism governing the dynamics of many diverse self-organized and adaptive systems [93, 94, 95, 96].

In chess—and sports in general—the balance of this trade-off may depend on skills. I thus investigate the extent to which skill level influences the approach to the game. In particular, I study the diversity in the player's arsenal of game openings across different Glicko-2 ratings. I calculate the Shannon entropy of the distribution (see *Methods*) of first move as white for each player and report the results in Fig 3.12. I find that beginners tend to open games with a diverse collection of first moves (as white) when compared to stronger players. Thus, the analysis captures beginners exploring a wider variety of first moves than experts, who instead are likely to begin with a typical move. At a first glance, this result might seem surprising, as skilled players are supposed to have better knowledge of opening theory. Yet, this may be linked to the ability of more skilled players to easily transpose into different opening variations in the following moves. Better awareness of transposition theory among experts



Figure 3.13. *Diversity in the black's response to white's first move.*

Boxplots showing diversity of black's response experienced by white player, for each of white's top 5 most played first moves- e4, d4, N f3, c4 and e3 (in descending order of popularity). As white, weakest players encounter lowest diversity in responses captured by low response entropy– for all of white's most played opening moves, except Nf3.

may allow them to reach many different openings from the same starting move, thus potentially eliminating the need to diversify in the first move itself.

So, overall, do experts specialize at the cost of diversity? To investigate further, I ask—how does skill level determine response diversity (as black)? For the top 5 white moves observed—*e*4, *d*4, *Nf*3, *c*4 and *e*3, I group the games of each player based on these moves and calculate the response diversity of the black to the white player. Results are shown in the different boxplots of Fig 3.13. Surprisingly, I observe a contrasting result. As white, beginners encounter the lowest diversity in black responses. This is captured by the low response entropy for all 5 of white's most played opening moves. Hence, beginners lack experience to the plethora of possible responses, which perhaps leaves gaps in their game.

Lastly, I point out that this increase in the diversity of responses at higher skill levels, might be what prevents players from increasing their Glicko-2 rating, as the potential to be surprised by your opponent keeps increasing as one climbs the skill ladder.

From the first move onward, players enter into established chess theory, where the many top variations of opening moves are well-explored. The next natural question to ask at this point is—How do players diversify beyond the first move as player move into opening theory? The beginning usually plays out like a well-choreographed dance, evolving in already classified opening sequences with standard names such as "Sicilian Defense", "Queen's Pawn Game", and so on. In Fig. 3.14, I show the top 9 openings used by players



Figure 3.14. Top 9 favorite openings among all players on lichess.

on *lichess.com*. Focusing on such opening sequences, I explore the specialization players achieve in the opening sequence. Results are shown in Fig. 3.15, where I define the "favorite opening" of a player as the most used one, assuming it is played at least 100 times.

Interestingly, the majority of players end up in their favorite openings only around 10% to 30% of the time. Furthermore, I find that expert players start with their favorite opening significantly more times than their second favourite. This is marked by the distribution falling below the diagonal line. Contrarily, beginners lie much closer to the diagonal, indicating that their favorite opening is played comparably to the runner up, thus pointing out a lack of specialization in a single opening.

Further analyses reveal that expert playing behavior comes in a variety of shapes and sizes, i.e., there are players who specialize and players who flexibly switch openings (diversify), see Fig. 3.15, column 4.

At the individual level, I find on average less diversity in opening selection (main lines) among experts, as shown in Fig. 3.16. As mentioned earlier, the ability to arrive into known openings through *transposition*, i.e., different sequences of moves that players may use to reach the same final configuration, might be unique to expert players. Arriving into fewer openings may allow experts to use learned chess theory and use optimal moves from memory, saving crucial time and preventing build-up of mental fatigue during the game.

However, accounting for the many different *variations* of the openings (see *Methods*), it is the experts instead who encounter the most diversity. This hints that experts like to enter into certain main openings—perhaps the ones they specialize in—which they follow-up by expanding their repertoire in the *vari*-



Figure 3.15. Usage of the top-2 favourite openings of a player.

Fraction of times players use their top-2 favourite openings. Different panels correspond to different skill levels. Density plots are used for better visualization. Expert players play more often their favorite opening sequence as compared to beginners.



Figure 3.16. Diversity and specialization in the opening sequence of moves is governed by skill level.

Distribution of diversity (entropy) of openings calculated for players of four different skill levels. Main lines and the variations are respectively depicted as solid and dashed curves.



Figure 3.17. Opening switches in a player's career properly normalised with the null model

Distribution of the number of opening switches in a player's career properly normalised with the null model aggregated into 4 skill categories. For the null model, I reshuffle the temporal order of the associated sequence of games, thus preserving the total number of victories, losses and draws.

ations to surprise opponents and catch them off-guard, a strategy not unique to chess but key in many competitive sports. Furthermore, upon investigating temporal organization of openings (main lines) used by a player, I find that experts switch openings between consecutive games more often than beginners (see Fig. 3.17). Thus, experts encounter higher temporal diversity in openings.

3.7 Performance and opening diversity

At this point one might wonder—how much exactly does specialized knowledge of favorite openings aid in victory? A naive argument would suggest that players would tend to prefer those main lines that give them the best results. If this is the case, the favourite opening of each player—the one mostly used would be the one that gives the best performance, that is the highest winrate. To investigate this, I calculate for each player the winrate of each of the player's top-3 most played openings and plot it against the frequency of their use. Results are shown in Fig. 3.18 for a sample of the players. Surprisingly, there are players whose top used opening performs worse than their lesser used openings. Besides, optimal players (black curves)—those who play more often their better performing openings—are just a few.

To quantify this effect in the whole population, I calculate for each player the difference in the winrate of the most played opening and the second most played one, showing its distribution in Fig. 3.19. The analysis reveals that when expert players do encounter their favorite opening, their winrate is more likely to be lower than their second favourite opening, when compared to beginners.



Figure 3.18. Winrate of top 3 openings of a player against opening frequency.

Each connected curve corresponds to a player. I show 15 random players who play at least 100 games with each of their top 3 openings. Curves of players whose winrate increases monotonically with the frequency of the associated opening are depicted in black and are deemed optimal.



Figure 3.19. Sub-optimal opening encounters.

Distribution of difference δw in winrate of associated to favourite and second favourite opening and winrate of their second favourite opening for the whole population of players. Different curves correspond to different skill levels. Dashed lines indicate mean values of the distributions. Stronger players encounter less optimal openings more often than weaker players. I note that players who do better in their second most played opening—as compared to their most played one—are experiencing sub-optimal opening encounters. Thus, I find that stronger players encounter sub-optimal openings more often than weaker players. In other sports, players are known to change their strategies over opponents and games so that specific habits are not exploited by opponents for strategic advantage. I speculate that such variations from optimal strategy might serve only a minor role in online play, where opponents are randomly selected from a large pool of millions of players. Thus, discovered sub-optimal encounters may be an opportunity for players to improve.

3.8 Diversity vs Career stage

Lastly, I explore diversity as a function of different stages of players' careers. Selecting players with at least 3000 games, I split them into 3 equal stages: early (0-1k), mid (1k-2k), and late career (2k-3k). For each play, I compute opening diversity in the different career stages and report it in Fig. 3.20. For both the opening move (top panel) and the opening sequence (main lines) (bottom panel), I find that players explore more in the initial stages of their careers, becoming more specialized in later stages, perhaps exploiting the knowledge of certain openings they have learned.

3.9 Discussion

This chapter extends the broader theme of success in careers by exploring chess as a unique context in which individual trajectories unfold, shaped by skill and patterns of specialization. In this chapter, I propose chess as a "natural laboratory" to investigate human behavior and performance [97, 98, 99, 100]. Chess offers a unique case as it lacks a stochastic component, allowing performance to be directly tied to skill, quantified here through the Glicko-2 rating. Analyzing nearly 1 million careers on lichess.org, I found patterns, including hot and cold streaks, reflecting bursts of victories and losses, seen in other domains such as science and business [22, 101, 23]. These findings suggest a possible universality to performance cycles across domains. Further analysis revealed that beginner players experience more frequent winning streaks, while streak length relates directly to skill level. Yet, irrespective of skill, players often face longer periods of repeated failure

Even just looking at simple patterns in the openings—thus neglecting the full complexity of game sequences—, I was able to characterize individual playing behavior across different career stages. In particular, expert players were



Figure 3.20. *Diversity in opening with career stage.*

Diversity in the opening move (top) and opening sequence (bottom) of moves. In later parts of one's career, diversity decreases, players prefer certain openings and specialize in them—playing them more often.

shown to behave differently from the very first move of the game, displaying a lower diversity in openings. Looking at chess as a process of interactions and reactions, I focused on the black's response to the white player's moves, finding that experts encounter the highest diversity from black. However, after accounting for different variations within the openings I discovered that experts were more diverse instead, hinting at a deeper understanding of the complexity of the different variations within the same line. Such findings corroborates some very recent ideas on opening similarity and complexity independently presented in Ref. [102], focusing on prediction of future openings and opening preparation.

Looking at individual careers over time, opening diversity was found to decreases at their later stages, pointing towards higher specialization as a player becomes more experienced. In addition, experts tend to play their favorite opening sequence much more than beginners, providing evidence for a tendency towards specialization. Nevertheless, counter-intuitively, I also found that players often do not have the ability to recognize their most successful opening, i.e. the one associated with the highest win-rate. Surprisingly, this is particularly true for more expert players, who have a higher chance of suboptimal encounters in opening, possibly because of the depth of responses and variations within opening lines coming from a skilled opponent. Indeed, in general decision making [103, 104, 105] and particularly in chess [99, 106], humans are known to work with heuristic approaches relying of intuition rather than searching for optimal solutions. The observed deviation from optimal strategies in the data might hence also hint at the existence of adaptive strategies with "fast and frugal" heuristics by players [107].

The analysis I have presented has some limitations. First, it focuses solely on openings—one of several chess phases. Nevertheless, this simple approach proved to be enough to reveal how experts differ from beginners in simple quantifiable ways. It also complements existing work on recall abilities of players for chess positions as a function of skill level [98]. A first natural extension in this direction would consist in analysing also other parts of the game, such as middle game and endings. A second limitation is that, when associating a skill level to a player, I inevitably considered Glicko-2 rating as a static, immutable measure. Instead, this rating systems is clearly in constant evolution throughout the career of a player. While including this dynamical aspect of ranking would surely add a missing aspect to the analysis, it is worth stressing that the measure is still a good proxy for skill level, as I have neglected the initial phase of the careers—associated to the steepest growth/change in Glicko-2 rating. Third, we observe in Fig. 3.2 that the number of games played by a player increases with skill level of the player. At the same time, the number of games also influences the opening behavior of the players as seen from career stage analysis in Fig. 3.20. Together this suggests that in future work, "number of games played" should be used as a control in addition to player ratings by which the players are currently stratified...

Taken together, this chapter represents a first step towards understanding the game mechanisms associated to performance in the careers of chess players. Future work might enrich this analysis by considering the complexity of chess games as a whole via considering the full sequences of moves instead of focusing on the important phases of the game only. Taken together, this chapter marks a step toward understanding performance mechanisms across competitive careers, with chess as a model. Future research could broaden this investigation to other sports, such as Go, tennis, or boxing, where opening moves and responses critically shape success and successful careers.

Code availability

The code used in this chapter is available at https://github.com/chowdharysandeep/lichess.git. 44 CHAPTER 3. QUANTIFYING PERFORMANCE IN CHESS CAREERS

CHAPTER 4

TEAM CAREERS IN SCIENCE: FORMATION, COMPOSITION AND SUCCESS OF PERSISTENT COLLABORATIONS

4.1 Introduction

Moving beyond success of individual entities, in this chapter, I focus on team careers. Teams are engines of innovation that propel collective scientific discovery. Large-scale collaborations such as CERN, the Human Genome project, LIGO, LHC, and the International Space Station are well-known endeavors created to expand the human frontier. Even at a smaller scale, the majority of research is increasingly produced by teams of pairs and triads of scientists, not individual researchers [48]. Moreover, recent work revealed that small teams tend to disrupt science and technology while large teams develop past research [56], highlighting the importance of underpinning the dynamics of smaller teams. In the age of the internet, research teams leave footprints in their lives that can be tracked retrospectively using data. The recent creation of large datasets documenting the evolution of science allows investigation of the multifaceted landscape of academia – from scientists and collaborations to funders, publishers, and institutions. Indeed these rapid developments have led to the emergence of the field of Science of science [16, 19], which aims to quantify the machinery underlying scientific production.

Nowadays much is known about the patterns in careers of individual scientists and what makes them impactful. For instance, scientific careers are known to display random occurrences of the most impactful work [21] and the presence of hot streaks [22, 23]. Role of mobility on impact [30], changes in research-interests [108] and an increasing trend to switch topics [31] have also been quantified in individual careers. Last decades however have seen a regime shift towards team science with teams increasingly dominating solo authors in knowledge creation, not only in sheer volume but also the attention they receive in the form of citations [47, 48]. The increasing predominance and success of teams have been linked to the increasing need for specialized knowledge and interdisciplinary collaborations in solving modern-day scientific problems [49, 50]. Inevitably, determiners of team impact are a highly active area of investigation. Factors such as– team size [48, 56] and diversity in geography and affiliation [50, 109, 110, 111], ethnicity [58], gender [112] and team freshness [113]— are known to significantly affect scientific impact. Mechanisms governing and analytical models of team assembly [47, 59] have also garnered some attention.

Although many factors governing scientific team assembly and determiners of impact have been studied, teams in science have only been viewed from the static lens. So far, team careers have received little attention. Few works that do consider the temporal aspect of collaborations are limited to ego-centric [62] or pairwise analysis [64]. Temporal dimensions of teams in science remains unexplored territory and the anatomy of a team's career remains largely unknown. Teams follow a dynamic lifecycle, akin to the trajectory of individual scientists' careers. They come together (form), achieve varying levels of productivity and impact, followed by eventual disbanding. Here I perform a large-scale study of teams careers in science. I extract the core members of a team, allowing me to track the life trajectory of a persistent team. Using a statistically-validated approach to disentangle the higher-order network of collaborations [114] I identify the scientific cores in teams and analyse half-a-million teams. I discern the patterns of formation, production, composition and impact of persistent scientific collaborations.

4.2 Methods

Extracting persistent cores

In this chapter I analysed publication data from OpenAlex [13, 14]. This database provides publication metadata, topic classification and citation records for around 205 million journal papers since 1900, covering 90 million scientists disambiguated using machine learning algorithms and integration with ORCID ids of scientists to identify authors [115]. I curate all scientists

with substantial publication records, keeping those with at least 20 papers, resulting in 4,000,926 scientists. In this context, identifying the core members of long-lasting scientific teams corresponds to detecting the maximal sets of significantly co-publishing authors [116]. In fact, these sets represent groups of authors that consistently work together, pruned from the members who only occasionally have published with them. To extract those, I start by constructing the underlying hypergraph of scientific collaborations, i.e. a generalized network that naturally encodes group relationships, usually called hyperedges[117]. Specifically, I construct a hypergraph whose nodes are authors and whose hyperedges represent joint publications among them.

In order to assess the statistical significance of collective interactions among sets of nodes in a hypergraph, one needs to take into account the heterogeneity of node activity. Indeed, if on one side the easiest way would be setting a fixed threshold on the minimum number of repeated interactions needed to consider a group of nodes a statistically validated set, this approach is sub-optimal, due to the multiscale nature of collaboration networks and the varying activity levels among scientists (the same threshold can be too restrictive for a author with limited publication records and very permissive for a more prolific author). To overcome this limitation, the method is based on a null hypothesis approach, that naturally tunes this threshold on the activity of the involved nodes. For simplicity, I start from the case of 3 nodes i, j, k that co-interact N_{ijk} times in hyperedges of size $n \ge 3$. The three nodes appear respectively in N_i, N_j, N_k hyperedges. Under the null hypothesis that each node selects randomly the hyperedges to which it participates - and thus its n - 1 counterparts in a hyperedge of size n - the probability of observing i, j, k interacting N_{ijk} times is

$$p(N_{ijk}) = \sum_{X} H(X|N, N_i, N_j) \times H(N_{ijk}|N, X, N_k)$$

= $\frac{1}{\binom{N}{N_i}\binom{N}{N_k}} \sum_{X} \binom{N_i}{X} \binom{N-N_i}{N_j-X} \binom{X}{N_{ijk}} \binom{N-X}{N_k-N_{ijk}},$ (4.1)

where $H(N_{AB}|N, N_A, N_B)$ is the hypergeometric distribution that computes the probability of having an intersection of size N_{AB} between two sets A and B of size N_A and N_B given N total elements. The probability $p(N_{ijk})$ in Eq. 4.1 represents the probability of having a random intersection of size N_{ijk} between the three sets of hyperedges of nodes i, j, k out of N total hyperedges [118], and is obtained through the convolution of two instances of the hypergeometric distribution. Starting from Eq. 4.1, I then compute a p-value for the triplet that

contains *i*,*j* and *k* through the survival function,

$$p(x \ge N_{ijk}) = 1 - \sum_{x=0}^{N_{ijk}-1} p(x).$$
 (4.2)

The p-value represents the probability of observing N_{ijk} or more hyperedges that contain - but are not limited to - the nodes i, j, k. The smaller the p-value, the higher the possibility that i, j, k constitute a significant set of size 3.

In the approach of [114, 116], N corresponds to the total number of papers in the corpus. In this case, this choice for the value of *N* implies that all hyperedges are equally accessible by all nodes. It is, however, unrealistic that a scientist could have participated in all papers, due to limited time and resources. Thus, I decided to bound the number of papers a scientist or a group of scientists could have worked on. Specifically, I approximate *N* by summing the number of unique publications by the scientists and all their coauthors, for each of the members of the group I am testing for significance.

For a generic hyperedge of *n* nodes, Eq. 4.1 becomes

$$p(N_{1...n}) = \sum_{X_{12}} H(X_{12}|N, N_1, N_2) \times \\ \times \sum_{X_{123}} H(X_{123}|N, X_{12}, N_3) \times ... \\ ... \times \sum_{X_{12...n-1}} H(X_{12...n-1}|N, X_{12...n-2}, N_{n-1}) \times \\ \times H(N_{12...n}|N, X_{12...n-1}, N_n).$$

$$(4.3)$$

How to set a rigorous criterion to assess whether a group of *n* scientists is statistically significant, once I have calculated the associated p-value? In order to do so, I test all p-values against a threshold of statistical significance α , after including a multiple hypothesis test correction which is needed because of the high number of tests - one per each group. In all the results presented in this chapter I use $\alpha = 0.01$. Coherently with the approach in [116]. I consider all smaller combinations of nodes constituting a significant set to be themselves significant sets. In other words, if the interaction *i*, *j* and *k* is significant, I do not test also the three couple obtained through combination of the triplet. This means that I start testing from the largest set and I then proceed towards the smallest. If a set of size n passesthe statistical test and is thus selected as significant because it rejects the null hypothesis, I do not test any of its smaller subsets. In this way, the obtained statistically significant sets can be considered maximal, i.e. for each of them there is no larger set that includes all its members that is also statistically significant.

The extraction of team cores resulted in over half-a-million persistent collaborations, with size ranging from 2 to 10. Yet, I limit the analysis to cores of size 2 to 6, as bigger cores are rare and statistics are insufficient for an in-depth analysis.

Knowledge broadness and knowledge diversity

In later analysis, I will use knowledge broadness and diversity to characterise team compositions. Here I show the method of evaluation for these measures. I start from the OpenAlex dataset, where researchers are associated with different scientific concepts with a score varying between 0 and 100. The concepts are organized in a hierarchy with 19 root-level concepts representing major disciplines such as physics, chemistry, computer science and so on, and 5 layers of descendants branching out from them, for a total of 65,000 concepts [119]. For this analysis, I consider the topmost layer of the concept hierarchy, namely the scientific disciplines. For each member of the team, I evaluate a knowledge vector where each component represents how strongly that scientist is associated with a scientific discipline. Each entry of the vector consists of the topic score normalized by the sum of the scores.

Knowledge broadness captures the breadth of the combined expertise of the persistent core. To calculate it, I first consider the sum of the knowledge vectors of the core members, normalized so that the components of the combined knowledge vector sum to 1. I then evaluate the entropy of this combined knowledge vector, calculated as $-\sum_{x} p_{x} log_{2} p_{x}$, where x is a scientific discipline and p_x represent how strongly the team is connected to it. Finally, to obtain the knowledge broadness, the entropy of the vector is normalized by its maximum theoretical value of $log_2(N)$ which corresponds to the case where the team is uniformly spread across all N disciplines. In particular, as Openalex has 19 disciplines in the top layer of its concept hierarchy here I have N = 19. For illustration, consider a core of three members, A, B, and C, whose knowledge vectors are A:[physics: 0.5, chemistry:0.5], B:[physics:0.7, chemistry: 0.1, biology: 0.2], C:[physics:1]. The combined knowledge vector for this team would be: [physics:0.85, chemistry: 0.077, biology: 0.077]. The knowledge broadness of the team is thus ~ 0.18 . The value of knowledge broadness ranges from 0 to 1. A value of 0 indicates a team where all members are associated to a single topic, i.e., a monodisciplinary team, while a value of 1 corresponds to a team in which the joint knowledge vector of the team is uniformly distributed across all possible topics in the data.

Knowledge diversity, on the other hand, quantifies how much team members are different with regard to the scientific disciplines they are associated with. Given the knowledge vectors of the core members, I quantify the knowledge diversity as 1 minus the the average cosine similarity between all pairs of knowledge vectors. For instance, in the example above, the cosine similarity between members A and B (~ 0.77) is calculated, then B and C (~ 0.95), then



Figure 4.1. *Identification of persistent core members in a team's career* A network of scientific collaborations (grey-shaded areas) is built based on the publication records. Then, two groups of scientists are identified in a team: core members (red), who consistently collaborate, and transient members (other colors), who publish together with core members only occasionally.

C and A (\sim 0.71). The values are then averaged and subtracted from 1 to obtain the knowledge diversity of the core (\sim 0.19). Knowledge diversity ranges between 0 and 1, where a value of 0 indicates all members have the identical knowledge vectors, i.e., all members are associated with the same strength to the same disciplines, whereas a value of 1 indicates that knowledge vectors of all members are non-overlapping, i.e., they are associated to completely different scientific disciplines.

Knowledge broadness and knowledge diversity cover complementary dimensions of team topic composition, as the first describes how different are the topics associated to the team, while the second captures how the members of the team are diverse among each other. For instance, a collaboration among scientists working in the same discipline will have low knowledge diversity and low knowledge broadness, while a team where members are associated with the same set of multiple disciplines will have low knowledge diversity but high knowledge broadness.

4.3 Formation, productivity and dissolution of persistent scientific collaborations.

Not all members of a scientific team collaborate in the same way. In science, a team is often composed of core members, who work persistently and recur-

rently together over the years, surrounded by transient members, who come and go in the collaboration. Since the exact set of co-authors can change from paper to paper, this makes following the trajectories of teams much more complicated than the careers of single individuals. To study team careers, my first step is thus to identify persistent scientific collaborations from empirical data. I analyze a large dataset collected from OpenAlex [13, 14], consisting of 248 million journal papers published since 1900, covering 90 million scientists across various scientific disciplines. From the publication records, I build a hypergraph [117] of scientific collaborations, where each hyperedge encodes the set of co-authors of a paper, and use a statistically-validated approach [114, 116] to extract those groups of scientists that have persistently published together over their careers (Fig. 4.1, see the Methods for a detailed description of the data and the methodology). Through this procedure, I identify 511,550 persistent scientific collaborations.

I begin by investigating the typical number of members in persistent scientific teams. I compute the distribution of core sizes (Fig. 4.2), finding that the greatest percentage of cores (42%) have 3 members, followed by cores of size 2 (31%) and 4 (20%). Larger persistent collaborations are rarer, with the fraction of cores of size 7 or higher being less than 1 in 100.

Typically, the formation of a scientific team requires a significant amount of time. Indeed, a persistent team may start as a small number of scientists working together, and gather further members around it later on. I thus ask: How fast do cores assemble? To examine this, I evaluate the time elapsing from the first publication authored by any subgroup of core members to the first publication authored by all members of the team. I refer to this quantity as the formation time. Note that, by my definition, teams of two members are established with a formation time of zero. The average formation time for cores of different sizes is shown in Fig. 4.3. I observe that smaller cores gather faster than bigger ones. In particular, teams of three members typically take 4.6 years to form, while larger cores take more and more time to gather, i.e., 7.3 years for cores of size 4, 9.2 years for cores of size 5, and 10.4 years for cores of size six. Furthermore, I note that a certain percentage of cores are formed instantaneously, i.e., they have a formation time of zero, meaning that the first publication by the team includes all of its core members. Out of all cores of size 3, nearly 16.8% assembled instantaneously. For larger cores, this number drops, 5% (for size 4), 2.3% (size 6), and only 1.4% (for size 6).

After their formation, teams start collaborating and publishing about their research. While the efficiency of scientific teams is hard to quantify, as I have no information on how long a team worked on a publication, I can analyze their career in terms of the number of joint publications as a function of time. To this



Figure 4.2. Distribution of core sizes.

Distribution of the number of scientists per team core. The greatest percentage of cores have 3 members.



Figure 4.3. Formation time of persistent teams.

Formation time of cores as a function of core size. Smaller cores take less time to form compared to larger ones.

end, I consider all cores of a given size that have published a certain number of scientific papers and measure the average time taken to produce them (Fig. 4.4). Moreover, I compute the yearly-production-rate, i.e., the productivity, for different core sizes (inset). The analysis reveals that bigger cores outproduce smaller cores, taking less time on average to publish the same number of scientific articles. Such a result calls for an in-depth analysis of how persistent scientific teams communicate, coordinate and organize as a function of the number of their members [55]. This observation may also partly explain the increasing dominance of teams in the production of science and arts [48].

Though they may persist for a long time, all scientific collaborations eventually end. Therefore, I examine the typical lifespans of persistent scientific teams. I calculate the survival probability of cores as a function of the career length, i.e., the time elapsed from their first to their last publication. I find that smaller cores are typically more persistent (Fig. 4.5). For instance, nearly 45% of



Figure 4.4. Productivity of persistent team cores.

Production time for a given number of papers for cores of different sizes. Inset: average number of papers published per year as a function of the core size. Bigger cores produce research articles at a faster rate.

scientist-duos work together for at least 5 years, while some of them can keep working together for as long as 30 years after their first joint publication. Larger cores, instead, have shorter careers: It is highly atypical for team cores of 4 or more scientists to continue publishing together after 10 years since their first publication. A possible explanation might be that a common grant limit the lifetimes of larger cores to 5-6 years (typical funding timelines). By combining our data with funding data, this question might be answered.

4.4 Composition.

The members of a persistent collaboration can be identified by various characteristics, including age, affiliations with universities, and scientific expertise. Understanding the composition of persistent teams in terms of the individual characteristics of their members, i.e., whether they overlap, match, or integrate, can shed light on the mechanisms that facilitate long-lasting collaborations. For



Figure 4.5. *Dissolution of persistent team cores. Survival probability of a core as a function of the career length, i.e., the time since its formation. Smaller cores have longer lifespans.*

instance, a persistent team may show age-homophily among its members or, instead, it may comprise young researchers and older experienced scientists, e.g., a long-standing collaboration between mentor and mentee. To understand the role played by age in persistent collaboration, I compute the career age of each scientist in the team (i.e., the time passed since the scientist's first publication), which I then use to characterize the age-composition of persistent teams at the time of core creation (i.e., the first joint publication).

I divide scientist into three age groups, namely young ("Y", career age less than 7 years), emerging ("Em", between 7 to 14 years) and established scientists ("Es", more than 14 years). I thus assign each core to one of 7 possible categories based on the age groups of the members. I distinguish cores where all scientists belong to the same age group (only young, only emerging, or only established), cores where two age groups are represented (young + emerging, young + established, or emerging + established), and cores with scientists from all three groups. The percentage of each category in the data is shown for various core sizes in Fig. 4.6. For dyadic cores, the most predominant age composition at time of assembly is a young and an established scientists working



Figure 4.6. *Age composition of persistent scientific cores. Percentage of each possible age composition of the core as a function of core size.*

together, a typical Student-Professor motif. This is followed in frequency by the young+emerging motif. For triadic cores, the mixed age composition i.e. Y+Em+Es likely corresponding to the typical PhD-PostDoc-PI core, starts to be a prominent age composition. Still, the dominant composition is of only Young + Established researchers. Beyond triads, for core size 4 and higher, mixed cores become increasingly common. I also analyze how the formation time and career length of a core vary as a function of its age composition, revealing shorter assembly times and longer careers for teams featuring younger members (Fig. 4.9a and 4.10a).

Next, I characterize how diverse persistent collaborations are in terms of academic affiliation. I will refer to this as the affiliation diversity of the cores. To quantify it, for each paper of the core, I consider the set of affiliations that the members have at the time of publication, and evaluate the minimum number of affiliations needed to represent all members of the core. For illustration, consider a core of 3 members, A, B, and C, having the affiliations A:[MIT, UCSD], B:[MIT, UCSD], C:[MIT]. In this case, the minimum number of affiliations needed to represent the core is 1, as all members share one affiliation (MIT). Instead, if I consider the case A:[Caltech, UCSD], B:[UCSD, Indiana], C:[Caltech], I would need at least 2 affiliations (Caltech and Indiana) to cover all members of the core. As the core members can change academic affiliation during the lifetime of the collaboration, I evaluate the minimum number of affiliations for each article published and define the affiliation diversity of the



Figure 4.7. *Co-localization of persistent scientific cores. Percentages of cores that are spread across 1 (mono), 2 (bi), 3 (tri), and 4 or more universities as a function of core size.*

core as the most common value. Fig. 4.7 shows the fraction of cores of different sizes being covered by one, two, or three universities. Possibly surprisingly and in contrast with an increasing trend to collaborate remotely [51], I note that 75.8% of cores are situated at the same university, hinting at the importance that common institutions play in sustaining long-term collaboration. It is also worth noting that cores that are not based at the same university are usually located at 2 universities (21.7%), while only very rarely persistent cores span 3 or more institutions (2.5%). Geographically speaking, 78.1% (84.7%) of persistent cores are situated within the same country (continent). I also find that co-presence at the same institution is associated with a shorter formation time and a longer lifespan of the core (Fig. 4.9b and 4.10b).

Finally, I aim to comprehend how diverse are core members in terms of their scientific expertise, namely whether the teams are mono-disciplinary or interdisciplinary. To achieve a data-driven understanding of the extent to which topic composition affects persistent collaborations, I measure two complementary dimensions of team interdisciplinarity, namely knowledge broadness and knowledge diversity. Knowledge broadness captures the breadth of the combined expertise of the persistent core. knowledge diversity, by contrast, quantifies how much team members are diverse with regards to the disciplines they are associated with. Both measures range from 0 to 1, where 1 are achieved for maximum values of multidisciplinarity of the team (broadness) and topic



Figure 4.8. Interdisciplinary diversity of persistent scientific cores.

Interdisciplinary diversity. The joint distribution knowledge diversity (quantifying crossmember conceptual distances) and knowledge broadness of the team (entropy of the sum of individual concept vectors of members). For monodisciplinary cores where all members belong to the same one field both knowledge broadness and core diversity are zero.

complementarity across members (diversity, see Methods for details).

Fig. 4.8 shows the joint distribution of knowledge broadness and knowledge diversity across the cores. I observe that only a small fraction of teams, almost 3%, are completely mono-disciplinary, namely all members work in one scientific field. Also, I notice that the density of cores tapers off as diversity increases, with few cores having a diversity larger than 0.5. This suggests that the core members should have a minimum amount of disciplinary overlap for their collaborations to persist. The vast majority of cores, however, are broad in their knowledge base (over 50% have more than 0.6). In particular, a significant proportion of cores ($\sim 9.1\%$) display broadness larger than 0.75, indicating teams which work at the interface of several disciplinary fields. Yet, knowledge

CHAPTER 4. TEAM CAREERS IN SCIENCE: FORMATION, COMPOSITION 58 AND SUCCESS OF PERSISTENT COLLABORATIONS

diversity never gets close to its maximum value, which would correspond to a team where no shared topics exist between its members. This highlights the importance of topic overlap to sustain persistent collaborations. Moreover, I find that teams with low knowledge diversity across members take less time to form and have longer careers, further supporting the association between team persistence and topic synergy, whereas low knowledge broadness, on the other hand, is associated with shorter lifespans and a slower formation process (Fig. 4.9c and 4.10c).



Figure 4.9. Formation time as a function of core knowledge composition.



Figure 4.10. Career length as a function of core compositions.

Core exclusivity as a function of core composition

As the members of a persistent core may work exclusively together or explore collaborations outside the core. I define the exclusivity of a persistent team as the ratio of the number of publications authored by all core members to the total
numbers of papers featuring at least one core member. In general, core tend to be highly exclusive, as more than 30% of all papers considered are published within a persistent collaboration. In terms of team composition, I do not observe a significant difference between mono and multi-university cores (Fig.4.11a). Instead, I note that cores made only of young researchers are in general more exclusive (Fig.4.11b). Also, I find that exclusivity decreases as a function of the knowledge broadness, while it is not strongly correlated with knowledge diversity (Fig.4.11c).



Figure 4.11. *Core exclusivity as a function of a) diversity in the academic affiliation, b) age composition, and c) knowledge broadness and knowledge diversity.*

Contribution of transient, non-core members of teams

A team is made by core members that persistently work together and transient members that sporadically publish with the core. So far, I have focused only on the core members of scientific teams. Here, I expand my analysis of teams to this latter group of members. First, I evaluate the fraction of papers published by persistent collaborations that include a given number of transient members (Fig.4.12a). I find that over 46.5% of papers are authored exclusively by core members, while 24.3% (14.0%) include 1 (2) non-core members, and less than 20% feature 3 or more transient members.

When categorizing cores by age composition and comparing this to the ages of non-core members (Fig.4.12b), I find that cores of young scientists tend to collaborate with younger transient researchers, whereas more established cores often work with older non-core members.

Incorporating non-core members into teams generally enhances knowledge diversity (for +87.7% of the core teams) and knowledge broadness (+76.7% of the cores), indicating that transient members add to persistent collaborations

CHAPTER 4. TEAM CAREERS IN SCIENCE: FORMATION, COMPOSITION 60 AND SUCCESS OF PERSISTENT COLLABORATIONS

a different expertise (Fig.4.12c,d). However, for 12.3% of cores, knowledge diversity decreases with the addition of a transient member, suggesting potential redundancy, which warrants further investigation. Finally, in terms of affiliations, in 95.0% cases non-core members are based in the same university or institute.



Figure 4.12. Contribution of transient members of persistent teams.

a) Percentage of papers published as a function of the additional number of non-core members.
b) Age distributions of the non-core members as a function of the age composition of the core.
c) Percentage change in the team diversity upon addition of non-core members.
d) Percentage change in the team diversity upon addition of non-core members.

4.5 Team success.



Figure 4.13. Example of a career of a scientific team.

Number of citations received by publications of a persistent collaboration involving the Nobel laureate Richard Henderson (Nobel Prize in Chemistry, 2017).

Peak performance in individual scientific careers is known to be randomly distributed, a result known as the random impact rule [21]. In other words, the most-cited article in a scientist's career can be, with an equal probability, any paper they published, from the first to the very last publication. Analyzing the publication history of persistent scientific collaborations (Fig. 4.13), I aim to understand the pattern of success in team career, and how their composition affects their impact. To examine whether team careers also exhibit the random impact rule, I adopt a similar methodology as [21, 120] and measure the relative position N* of the highest-impact paper in a core's career, i.e., in the sequence of its N publications. I measure the impact as the number of citations after five years (c5), normalized to account for inflation and large variations between different disciplines [121]. Fig. 4.14 shows the cumulative distribution function $P(\leq N^*/N)$, namely the probability that the most-cited work of a core appears before the N*-th publication.



Figure 4.14. Temporal dynamics of team success.

Cumulative distribution $P(\leq N^*/N)$ *for cores with joint total papers* $(N) \geq 10$ *, where* N^*/N *denotes the order of the highest-impact paper in a core's career, varying between 1/N and 1. The cumulative distribution of* N^*/N *is almost a straight line with slope* ≈ 1 *, indicating that* N^* *has roughly the same probability to occur anywhere in the sequence of papers published by a core.*

I observe that the function nearly follows the cumulative probability of a uniform distribution, which would indicate that the highest-impact article of a

CHAPTER 4. TEAM CAREERS IN SCIENCE: FORMATION, COMPOSITION 62 AND SUCCESS OF PERSISTENT COLLABORATIONS

core can be published at any point in the career. While the most-cited article can occur at any time, the average impact of a team throughout its career can show non-random patterns. In the evolution of a team's career, freshness decreases. This can have consequences for the success of the team and the impact it gathers. To test this, I split the papers of each team (ordered chronologically) into two halves and compute their average impact, measured as the normalized number of citations after five years from the publication. My analysis captures a higher average impact in the first half of the career compared to the second half, hinting that freshness in an early career may bring a higher impact (Fig.4.15). My findings are in agreement with results showing how the team freshness has a positive effect on the team impact [113].



Figure 4.15. Team impact as a function of career stage.

Boxplots show the distribution of the the average impact (c5) of the core papers in the first and second half of their careers.

Another universal feature of individual creative careers is the presence of hot-streaks, periods of clustered high-performance works observed in various contexts, from science [22, 23] to arts and other creative domains[122]. My results show a tendency for highest-impact publications to be clustered, i.e., to occur close to each other in the team's career, thus confirming the presence of hot-streaks ((Fig.4.16)).

Having characterized the diversity of persistent collaborations in terms of age, affiliations, and scientific expertise, I am now interested in understanding how those features contribute to their academic success. First, I measure the average paper impact, i.e., the average normalized number of citations after



Figure 4.16. Hot-streaks in team careers.

(a) Cumulative distributions $P(\leq (N^*-N^{**})/N)$, where N^*-N^{**} denotes the distances in order of the highest and the second-highest impact papers in a core's career respectively. (b) $P(\leq (N^{**-}N^{***})/N)$ where N^*-N^{***} where N^*-N^{***} denotes the distances in order of the highest and third-highest impact papers. (c) $P(\leq (N^*-N^{***})/N)$ where $N^{**}-N^{***}$ denotes the distances in order of the second-highest and third-highest impact papers. The shuffled model randomizes the order of impact within a teams career, in all three panels. Real data shows higher likelihood than shuffled model for highest impact work to occurs in bursts (close to each other in time), indicating the presence of hot-streaks in team careers.

five years, for all possible core age compositions in the data. Cores composed of members from the same age groups perform typically worse than collaborations with scientists from multiple age groups (Fig.4.17).



Figure 4.17. Team age composition and impact.

Cores are separated based on age composition of members and the distribution of the average impact per paper for each age composition is shown.

Besides the age composition, the working location of a core team can affect the impact of the work it produces. To test this, I compare the average paper impact after five years for mono and multi-university cores. In agreement with previous research [50], at the aggregated level I observe that multi-university teams are more successful than mono-university ones. However, a time-resolved analysis of the average impact, consisting in comparing mono

and multi-university cores as a function of the year of formation, reveals a more nuanced picture Fig. 4.18. I find that the impact advantage of multi-university core is a recent phenomenon. Mono-university cores formed before the 2000s have on average a higher impact compared to multi-university collaborations formed during the same years. However, while more recently formed monouniversity cores do not show a significantly higher impact compared to older collaborations, multi-university cores formed in the last decade have almost three times the impact of cores started in the 60s. Hence, multi-university cores formed after year 2000 are more successful than mono-university cores formed in the same period. A possible explanation for this observation lies in recent technological progress, from the development of computer-based mailing systems to the advent of internet and other communication technologies, which have enhanced the experience of remote collaboration, limiting the logistic advantage associated with working in close physical proximity.



Figure 4.18. *Team impact for mono-university and multi-university cores. The average impact of cores that formed in a given decade, where cores are separated into those co-located at the same university, and others with multiple institutions.*

In addition to impact, I look at the productivity of teams as a function of knowledge broadness, grouping the cores based on the level of knowledge diversity, i.e., low, medium, and high diversity, respectively (Fig. 4.19). I find that the relationship between productivity and knowledge broadness approximately follows a U shape.

This suggests that teams with a focused approach are more productive than those with moderate levels of diversity. However, productivity increases once



Figure 4.19. Team productivity as a function of knowledge composition of members.

The average impact of cores that formed in a given decade, where cores are separated into those co-located at the same university, and others with multiple institutions.

again for teams exhibiting highest degrees of broadness. For a given broadness, teams high in knowledge diversity show higher productivity. As an additional analysis, I study how productivity depends on the core composition in terms of age and academic affiliation, finding no significant impact of these features (Fig.4.20).

Next, I investigate team success as a function of the knowledge diversity and the knowledge broadness of its core members. In Fig. 4.21, I show the average c5 of the team publications as a function of the knowledge broadness, grouping again cores based on their knowledge diversity. I observe an inverted U shape relationship between impact and knowledge broadness of the core, for all classes of knowledge diversity, with intermediate values of broadness supporting highest impact. When the knowledge broadness of a core is very low, the impact of a team's work might be limited to narrow disciplinary fields. Yet, when a core's knowledge base becomes too broad, a team's work yields a lower impact. Similar observations were made at the level of single works, where highest impact is obtained in papers which display a balance between conventional and atypical combinations of prior work [52]. Furthermore, I observe that, for almost any value of knowledge broadness, knowledge diversity has a negative relationship with the average core impact. All in all, my findings complement previous results on the complex relation between disciplinary diversity



Figure 4.20. Productivity as a function of core's geographical and age composition.

and impact in teams[58].

I conclude this analysis by assessing the role of transient members on team performance. I classify core publications into two groups, namely those authored exclusively by core members and those including other contributors, and evaluate the average impact for these two categories. Only teams with at least one publication with non-core members and one with only core members are kept for a fair comparison. Fig. 4.22 shows the distribution of the average c5 across scientific cores for the two groups of papers. The analysis reveals that publications involving transient members generally show lower impact compared to those authored by the core only, suggesting that even though transient members add diversity to the team, they might not boost impact.

4.6 Discussion

In this chapter, I moved beyond success of individual careers and introduced the notion of team careers to unravel the determinants and the temporal patterns behind persistent collaborations in science. I investigated the features of half a million persistent teams and the patterns governing their formation and lifespans, composition, production and eventually impact.

Core teams of three scientists were prevalent, highlighting the need to study team careers beyond pairwise collaborations [64]. Larger cores formed slower, and had shorter career lengths. Persistent collaborations with smaller cores often featured a mix of young and established researchers, while larger cores included members from all age groups. Additionally, members of persistent cores were affiliated with the same university for the majority of cores. A wide range



Figure 4.21. *Success as a function of interdisciplinary diversity of the team. Average impact as a function of knowledge broadness of the core. Cores were also assigned categories according to core diversity into 3 equal-sized categories –low, medium or high.*

of diverse degrees of disciplinary composition was observed, including teams with large knowledge broadness, likely linked to the rising popularity of interdisciplinary research topics [123, 124]. Teams with a higher knowledge diversity were found to have shorter lifespans and longer formation times, highlighting the importance of topic synergy for persistence.

In my examination of temporal patterns in the success of persistent teams, I found that the highest-impact paper can occur at anytime, with the same probability, throughout the teams career. Over extended periods, however, I found that the average magnitude of impact is higher for publications in the first half of the team career, in agreement with earlier observations that freshness is associated with high multidisciplinary impact [113]. Besides, hot-streaks were observed in team careers, with the highest-impact papers showing a tendency to be clustered in time, mirroring findings on individual scientific [22, 23] as well as artistic careers [122].

Previous research shows that a time-aggregated analysis captures multiuniversity teams to be more successful [50], yet, I find that, for persistent collaborations, this is the case only since the 2000s. Coinciding with the advent of ICT technologies and the internet, remote research teams achieve higher im-



Figure 4.22. Effect of non-core members on team success. Distribution of average impact of team publications authored solely by the core members (blue) or involving transient, non-core members (green).

pact. Yet, recent studies indicate that remote work results in fewer bridges and encourages asynchronous communication among employees [125], while also being linked to a lower likelihood of disruptive scientific ideas as compared to onsite collaborations [126]. Impactful teams seemed to strike a balance in terms of knowledge broadness, while diversity among team members was associated to low citation accumulation, in accordance with previous analysis at the level of single manuscripts [52], where the highest impact is obtained for papers that display a balance between the conventional and atypical combinations of prior work. Contrasting with citation impact, productivity was highest for teams that were topic focused or ones that maximized knowledge breadth. Moreover, publications authored solely by core members yield higher impact compared to those involving transient non-core members, hinting that expanding teams may not always enhance outcomes for persistent collaborations.

While my analysis reveals features of successful persistent collaborations, whether the observed patterns reflect a causal relationship between team composition and success, or success is simply a prerequisite for collaborations to persist and survive, is an open question which might be clarified in future investigations. Besides, future works investigating team composition might benefit from integrating additional information about team members, including gender and ethnicity. Indeed, I know at the level of teams defined as single set of coauthors that gender diversity among members plays a role in the impact achieved by the publication [127], paving the way for a similar investigation on the role of gender diversity for persistent collaborations. I note that bibliometric databases do not include self-declared gender or race by authors and inferring such metadata from names often introduces algorithmic biases [128] which are distributed unevenly among other demographic traits [129], highlighting the difficulties in properly carrying out such analysis.

In the future, these results can be combined with prior theoretical work [47, 130, 59] to build a more complete theory of team assembly in science. Besides extracting persistent collaborations active over time, I can study team dynamics, how teams evolve and transform themselves and collectively adapt to a scientific ecosystem that constantly evolves in time [131]. Furthermore, a crossdisciplinary analysis comparing teams in rapidly evolving fields like AI with those in more stable areas like mathematics, can guide tailored team formation strategies.

Examining the role of funding on team careers can highlight how external support mechanisms can promote meaningful persistent collaborations. Indeed, the scientific ecosystem is largely shaped in response to funding. Over time, disciplines differ in terms of funding support they receive, with interdisciplinary research so far achieving lower funding success [132], likely impacting the longevity of such teams. Moreover, grant success is strongly dependent to the collaboration network [133, 79], particularly for young researchers, and this may incentivize participation in teams on strategic rather than scientific grounds. From the perspective of individual careers, how persistent collaborations affect they way in which researchers navigate the knowledge landscape is also an intriguing question to address [134]. Long-term career outcomes of individuals who participate in high-impact, persistent collaborative experiences can be compared to those who work in more transient or less successful teams. Notably, early-career collaborations with elite scientists predict subsequent career success [135]. Investigating this aspect further could provide insights into how early career researchers can strategically navigate collaborations and switch between teams to enhance their professional growth and impact.

Overall, my work identifies persistent teams in science and captures temporal as well as compositional patterns of success in shared careers. This research informs both scientists in building their collaborations and funders in choosing which research teams to support and promote.

Data availability

The data from openalex.org used in this chapter is openly accessible for down-load using the API https://api.openalex.org/works.

Code availability

The code used in this chapter is available at https://github.com/chowdhary-sandeep/sciscicareers.

CHAPTER 5

MODEL OF SOCIAL CONTAGION IN A POPULATION WITH EVOLVING GROUP STRUCTURE

5.1 Introduction

In the previous chapter, I explored team success in science, examining the roles of persistent core members and transient collaborators on the team performance. This chapter extends that exploration and investigates how these teams (groups of scientists) contribute to propagation of ideas throughout the scientific ecosystem, in particular, focusing on the role of temporal persistence of these groups. This line of research connects directly to broader questions about how social contagions spread across networks, a phenomenon extensively studied in fields from epidemiology to sociology [65]. From the viral spread of the Ice Bucket Challenge to the rapid adoption of smartphone technology globally, social norms and ideas cascade through networks, profoundly shaping collective behavior. Such contagion processes, including the spread of diseases, opinions, and rumors, are ubiquitous in nature [136, 137, 138]. In all such cases, the contact structure of the underlying population has a crucial role in determining the emerging collective behavior, making network science one of the primary tools to investigate spreading dynamics in real-world systems [139, 140, 65, 141, 142]. For instance, pioneering investigations have shown that heavy-tailed degree distributions in the contact structure lead to a vanishing epidemic threshold, a behavior which can not be observed neither in well-mixed population nor in homogeneous networks [143]. For the biological spread of pathogens, contagion is typically mediated by pairwise interactions,

where each link represents an independent source of infection. However, this mechanism of *simple* contagion does not seem to accurately describe social contagion. To acquire new ideas, norms or opinions, spreading is better modelled by *complex* contagion [67, 144, 145, 146]. In this case, individuals are subject to the simultaneous pressure of their neighbors, leading to a dynamics of cascades which has also been empirically observed in a number of different contexts [68, 69, 70, 71, 72].

For many years, the wide majority of networked systems have been represented by graphs, collection of edges and links, where interactions are naturally limited to dyadic ones [147, 148]. However, in most real-world networks, interactions can also occur among groups composed by three or more individuals. All these systems are better described by simplicial complexes or hypergraphs, which naturally take into account the presence of higher-order interactions, providing a suitable extension of the traditional network framework beyond pairwise interactions [149, 117, 150, 151]. In particular, simplicial contagion is a newly proposed paradigm that allows one to model at the microscopic scale the effect of group interactions (described as simplices of different order) on spreading dynamics [73]. Interestingly, if the infection rate associated to the higher-order interactions is high enough, this leads to the emergence of new collective behavior, making the transition from the healthy to the endemic phase explosive, and giving rise to metastable states. I point out that, while explosive phenomena are in general unusual in traditional epidemic processes [152], instances of such transitions have been observed in specific cases. A pertinent example is the one of cooperative [153] or synergistic contagion in networks [154], where a dynamical enhancement in spreading leads to an abrupt epidemic transition. Explosive transitions have also been observed in multiplex networks where the spreading dynamics in a layer is coupled to dynamical processes taking place on other layers [155, 156].

In context of higher-order interactions, such result was obtained analytically by a mean-field analysis and confirmed by numerical simulations [73, 157], has also been replicated under different modeling frameworks, such as the microscopic Markov chain approach [158], the generalised link equation [159], approximate master equations [160], and on different higher-order representations, such as hypergraphs [161, 162, 163]. The disruptive presence of higher-order interactions is not limited to contagion dynamics, as new collective behavior has also been observed in the case of synchronization phenomena [164, 165, 166, 167], random walk [168, 169], consensus [170, 171], ecological [172, 173] and evolutionary dynamics [174] when extended beyond simple dyadic ties. For pairwise contagion, the temporal nature of interactions, where links can be created and destroyed over time, is known to significantly affect the evolution and the long-term properties of the spreading process [175, 176]. Indeed, temporal networks [177] are routinely used as a modeling framework to properly capture diffusive processes taking place on realistic populations where the contact structure changes over time [178, 179, 180, 181]. Recently, also higher-order social networks have been found to have a non-trivial temporal dynamics [182]. Yet, so far very little attention has been devoted to understanding how temporality affects spreading on higher-order structures [183].

In this chapter, I extend models of simplicial contagion to the case of time-varying networks, where both pairwise and higher-order interactions can evolve over time. I compare the contagion process on static and temporal simplicial complexes. The dynamics of the static case presents bistability, meaning that the long-term behavior of the system is determined by the size of the initial seed of infectious nodes. I numerically characterize the basins of attraction of healthy and endemic states in static and temporal higher-order structures, showing that persistent temporal interactions anticipate the onset of the endemic state in finite-size systems. This means that the same number of initially infected agents might or might not lead to an endemic stationary state, depending on the temporal properties of the underlying network structure. To this aim, I propose a simple model to tune the degree of temporal correlations (persistence of groups) in synthetic structures that evolve over time, and investigate how this variable affects the long-term outcome of the spreading dynamics. I show that temporality can significantly reduce the enhancement of epidemics typically induced by higher-order contagion terms in the forward transition to the endemic state. By contrast, the backward transition to the infection-free state remains unaffected by presence of temporal correlation or lack thereof. Finally, I study simplicial contagion on temporal higher-order networks that present degree heterogeneity, showing once again that temporality hinders higher-order spreading, but in a less pronounced way than for homogeneous structures.

5.2 Model of contagion

I study social contagion in simplicial complexes which evolve over time. In particular, following Ref.[73], I consider an SIS model, where each one of the N interacting nodes can be in either of two states – susceptible (S) or infected (I). I consider interactions up to groups of three, such that 1-simplices (links) encode standard pairwise interactions, while 2-simplices describe three individuals interacting together (and this is structurally different from having three links that form a triangle). This choice to only consider 1 and 2-simplices (pairs and triads of nodes) is motivated by our observation in Chapter 4 featuring teams, where we discovered that the most common type of team cores are of size 2 and 3.



Figure 5.1. Temporal higher-order networks.

Schematic of a time-varying higher-order network where both pairwise and higher-order interactions evolve over time.

In a time-step of the SIS model, any infected individual can infect their susceptible neighbours connected by 1-simplices with a probability β_{\parallel} , and infected nodes can recover with probability μ and become susceptible again. However, in the simplicial version of the model, 2-simplices provide an additional way for a contagion event to happen. In particular, if a susceptible individual is part of a 2-simplex while the other two members of the simplex are infected, there is an additional probability β_{\triangle} to also get infected – associated to a microscopic description of social reinforcement induced by group interactions.

I write the discrete time evolution equation for the infection probabilities of each node at a particular instant using the Microscopic Markov Chain Approach (MMCA) [184]. MMCAs have been extended to temporal networks, allowing for an analytical computation of the epidemic threshold [179], and more recently to simplicial complexes, though in this context the non-linear term associated to contagion in 2-simplices only allows a numerical solution [158]. According to this approach, the probability of a generic node *i* to be infected at time t + 1 is

$$p_i(t + 1) = (1 - q_i(t)q_{i,\triangle}(t))(1 - p_i(t)) + (1 - \mu)p_i(t),$$
 (5.1)

where the first term on the right-hand side of Eq. (5.1) represents the probability at time *t* for a susceptible node to get infected. This is given by the product of $(1 - p_i(t))$, the probability that node *i* is susceptible, and $(1 - q_i(t)q_{i,\triangle}(t))$, the probability that *i* is infected by at least one of its neighbours. The second term, $(1 - \mu)p_i(t)$, stands for the probability that node *i* is already infected at time *t* and does not recover. Here $q_i(t)$ defines the probability that node *i* is not infected via pairwise interactions with its neighbours,

$$q_i(t) = \prod_{j \in \Gamma_i(t)} \left(1 - \beta_{|} p_j(t) \right) , \qquad (5.2)$$

with $\Gamma_i(t)$ denoting the set of 1-simplices containing node *i* at time *t*. Similarly, $q_{i,\triangle}(t)$ defines the probability that node *i* is not infected by any of its 2-simplicial interactions,

$$q_{i,\triangle}(t) = \prod_{j,\ell\in\triangle_i(t)} \left(1 - \beta_\triangle p_j(t) p_\ell(t)\right) , \qquad (5.3)$$

with $\triangle_i(t)$ denoting the set of 2-simplices containing node *i* at time *t*.

Notice how, in contrast with Ref. [158], here $\Gamma_i(t)$ and $\Delta_i(t)$ are functions of time, and allow us to generalize the MMCA approach to evolving simplicial complexes.

5.3 Social contagion on static and temporal simplicial complexes

I begin by comparing contagion processes in static simplicial complexes and in higher-order networks that change over time. A schematic of a time-varying higher-order network is shown in Fig.5.1 where 1-simplices and 2-simplices are respectively coloured in blue and yellow. As stated earlier, my focus on 1- and 2-simplices (pairs and triads of nodes) is motivated by our observation in Chapter 4 featuring teams, where we found that the most common type of team cores are of size 2 and 3.

In particular, I consider random simplicial complexes (RSCs) with N = 500 nodes generated following the algorithm introduced in Ref. [73]. While RSCs do not represent real dynamical systems, they provide a special (extreme) case with zero-temporal correlation for comparison against the static case where temporal correlation is maximum– a constraint we will relax later. The procedure allows to obtain homogeneous simplicial complexes with controlled generalised degree properties [185], namely $\langle k_{|} \rangle$, the standard pairwise degree, and $\langle k_{\Delta} \rangle$, the average number of 2-simplices incident on a node. In such model, 1-simplices are created akin to the Erdös-Rényi model, by connecting any pair (i, j) of vertices with probability $p_{|}$. Similarly, 2-simplices are added by connecting any triplet (i, j, ℓ) of vertices with probability p_{Δ} . For two desired values of $\langle k_{|} \rangle$ and $\langle k_{\Delta} \rangle$ it is possible to choose $p_{|}$ and p_{Δ} according to: $p = \frac{\langle k_{|} \rangle - 2\langle k_{\Delta} \rangle}{N-1-2\langle k_{\Delta} \rangle}$ and $p_{\Delta} = \frac{2\langle k_{\Delta} \rangle}{(N-1)(N-2)}$ [73].





I show the fraction of infected nodes at the equilibrium starting from a single infected node as a function of rescaled pairwise $\lambda_{|}$ and simplicial λ_{\triangle} infection rates for static (a) and temporal (b) simplicial complexes with N = 500 nodes. In the static case, the epidemic onset (solid black line) as a function of $\lambda_{|}$ is anticipated as I increase λ_{\triangle} . This suggests that the chosen initial infection of size $\frac{1}{N}$ belongs to the basin of the infection-free state for small values of λ_{\triangle} , moving into the basin of the endemic state upon increasing λ_{\triangle} . For time-evolving higher-order networks such effect is not observed, and I find a suppression of the endemic phase which can not be reached for low values of $\lambda_{|}$, independently on the value of λ_{\triangle} . The backward transition to the infection-free state (dashed black lines) is largely unaffected by the temporality of the interactions. I set $\mu = 0.1$, $\langle k_{|} \rangle = 12$ and $\langle k_{\triangle} \rangle = 5$ for both scenarios.



Figure 5.3. Size of the infected population as function of λ_{\parallel} obtained analytically in the mean-field limit for different values of λ_{\triangle} .

The basin of the infection-free state shrinks as λ_{\triangle} *is increased allowing earlier onset of endemic phase for the forward transition.*

I am particularly interested in studying how temporality affects the basins of attraction in the bistable regime which separate the endemic state from the infection-free state. Thus, I simulate the contagion process by first infecting a single node chosen at random and check whether this is sufficient or not to fall into the absorbing state with no epidemics. In particular, I numerically track the temporal evolution of the system at each time step *t* by updating the infection probabilities $p_i(t)$ for all nodes as dictated by Eq. (5.1). I iterate Eq. (5.1) for long time (10000 time steps) and compute the density of infected node in the stationary state by averaging the infection probabilities as $\rho = \frac{\sum_i p_i}{N}$.

In Fig. 5.2a I show ρ for a static RSC as a function of rescaled pairwise, $\lambda_{\parallel} = \beta_{\parallel} \frac{\langle k_{\parallel} \rangle}{\mu}$, and simplicial, $\lambda_{\triangle} = \beta_{\triangle} \frac{\langle k_{\triangle} \rangle}{\mu}$ infection parameters. In Fig. 5.2b I compute ρ for RSCs that change over time, where at each time *t* I generate a new realisation of the RSC model with the same $\langle k_{\parallel} \rangle$ and $\langle k_{\triangle} \rangle$ of the static simulations. In both heatmaps, two distinct regions separated by the black solid curves appear, an infection-free region where $\rho = 0$ and an endemic region where a macroscopic fraction of the nodes is infected.

In the static case, as I increase λ_{\triangle} , the epidemic onset occurs for progressively smaller values of $\lambda_{|}$ in finite-size systems. This means that the seed of infectious nodes of fixed size $\frac{1}{N}$ belongs to the basin of attraction of the infection-free state for small values of λ_{\triangle} , while it moves to the basin of the endemic state upon increasing λ_{\triangle} . Coherently with the results obtained with the mean-field formalism [73], above a critical value of $\lambda_{|}$, the system always reaches a non-zero fraction of infected agents which grows together with λ_{\triangle} . It is worth mentioning that in static structures [Fig. 5.2a] I find a slight anticipation of the



Figure 5.4. Effect of initial infection size on the onset of the endemic state.

(a) Density of infected nodes for static (b) and temporal (c) simplicial complexes as function of $\lambda_{|}$ for three different initial infections, $p_0 = \frac{0.1}{N}, \frac{0.5}{N}, \frac{1}{N}$ and two different values of rescaled simplicial infectivity, $\lambda_{\triangle} = 15$ (dashed curves) and $\lambda_{\triangle} = 30$ (solid curves). An early onset of the endemic phase is observed for sufficiently high values of the infected seeds and λ_{\triangle} with a MMCA approach, compatible with my observations in (a). By contrast, in temporal simplicial complexes, even for higher values of initial infection $\frac{1}{N}$ and high simplicial infectivity, (e.g. $\lambda_{\triangle} = 30$), there is a striking suppression of contagion and early onset of endemic state does not occur. I set $\mu = 0.1$, $\langle k_{|} \rangle = 12$ and $\langle k_{\triangle} \rangle = 5$ for both static and temporal scenarios.

epidemic threshold due to the MMCA as compared to the mean-field treatment, according to which the critical threshold $\lambda_{\parallel}^{c} = 1$ for $\lambda_{\triangle} = 0$. This is consistent with what has been already observed in Refs. [158, 159]. More interestingly, below this critical value, it is still possible to end up in the endemic state due to the higher-order contributions, but only if the seed of infectious nodes is big enough (critical mass). In this case, the system undergoes an abrupt transition.

Surprisingly, by contrast, λ_{\triangle} does not affect the onset of the epidemics in temporal simplicial complexes of finite size. This is clear from Fig. 5.2b, where the transition from the healthy to the endemic state is only observed as a function of $\lambda_{|}$, with the critical point $\lambda_{|}^{c} = 1$ coinciding with what predicted by the mean-field approach [73]. Notice indeed that critical mass effects are completely suppressed, and below $\lambda_{|}^{c}$ the same seed of infectious nodes can never sustain the epidemics –as opposed to what happens in the static case for sufficiently high values of λ_{Δ} .

So far I have focused on forward transitions from the infection-free state to the endemic state. Yet, abrupt transitions are typically associated to the emergence of hysteresis cycles. For this reason I also explore the backwards transition from the endemic phase to the infection-free state by choosing the stationary-state infection probabilities obtained at the higher value of λ_{\parallel} as the initial seeds for simulations at lower $\lambda_{|}$ values. I show the backward transitions as dashed black lines in Fig. 5.2a, 5.2b and find that they remain unaffected by temporality.

In Fig. 5.2, I fixed the size of the initial seed of infectious nodes at $\rho(0) = \frac{1}{N}$. To better characterize the two basins of attractions in the bistable regime and the associated critical mass effects, in the following analyses, I vary the initial seed size and numerically investigate the onset of the epidemic. In particular, in Fig. 5.3 I first show the analytical solution for the stationary ρ in the meanfield approximation derived in Ref. [73] as function of λ_{\perp} for different values of λ_{Δ} . The dashed curves represent the unstable solutions that separate the basin of the infection-free state ($\rho = 0$) from the endemic state ($\rho > 0$). As λ_{\triangle} is increased, I see that the basin of the infection-free state shrinks so that the endemic phase can be reached for progressively smaller values of initial infection size $\rho(0)$. Indeed, consistent with this, my numerical investigations on static simplicial complexes (Fig. 5.4a) reveal that while a small initial infection of size $p_0 = \frac{0.1}{N}$ does not lead to early onset of endemic phase no matter the value of λ_{\triangle} , increasing the initial seed size to $\frac{0.5}{N}$ or $\frac{1}{N}$ leads to early onsets on the endemic phase in the system with N = 500 nodes. As expected, the onset occurs even earlier for higher values of λ_{\wedge} . By contrast, in temporal simplicial complexes, as shown in Fig. 5.4b, the onset of the endemic phase in temporal simplicial complexes is largely independent of λ_{Δ} , consistently with what was observed in Fig. 5.2b. This suggests that the basin of the infection-free state shrinks fast in static simplicial complexes as λ_{\wedge} increases. As a consequence, the relevance of group effects is strongly mitigated when I consider temporality, a realistic feature of many real-world social systems.

5.4 Contagion on temporally correlated higherorder networks

In the previous section I observed that introducing time-evolving structures can significantly impact contagion on higher-order networks, by altering the basin of the infection-free state in finite-size simplicial complexes. However, the way in which network structures evolve can be different. For instance, a social system may change more or less quickly, giving rise to different temporal correlations among networks at consecutive times. I thus consider as a measure of temporal correlation:

$$\sigma = \frac{1}{2T} \sum_{t=1}^{T} \frac{n(|_t \cap |_{t+1})}{n(|_t \cup |_{t+1})} + \frac{n(\triangle_t \cap \triangle_{t+1})}{n(\triangle_t \cup \triangle_{t+1})}$$
(5.4)



Figure 5.5. A schematic of temporal simplicial complexes with low and high temporal correlations.

where Δ_t is the set of 2-simplices at time t and $|_t$ is the set of 1-simplices which are not part of any 2-simplex at time *t*, $n(\triangle_t \cap \triangle_{t+1})$ is the number of 2-simplices that persist from time t to the next time step t + 1 and $n(\Delta_t \cup \Delta_{t+1})$ is the total number of 2-simplices present at time t or t + 1. Analogously, $n(|_t \cap |_{t+1})$ and $n(|_t \cup |_{t+1})$ are defined for 1-simplices.

In order to investigate how the evolution of the network affects the spread of contagion, I introduce a model to systematically tune temporal correlations in simplicial complexes, where at each time the network is described by a RSC. In details, I recursively generate a new simplicial complex at time t + 1 by randomly rewiring with probability $f \in [0,1]$ the 1-simplices and 2-simplices present at time t. In this way, I am able to generate a temporal sequence of RSCs. Using such a model for sparse graphs, I can tune the temporal correlation σ in an effective range between 0, describing the absence of correlation, and 1, where network structure does not change over time. Two schematics of temporal simplicial complexes with low and high correlation are shown in Fig. 5.5.

In the following analysis, I focus on the forward transition to endemic state only, as the backward transition is unaffected by temporality as observed in Fig. 5.2 (dashed curves). I first infect a single node and simulate the epidemic process on top of two distinct sequences of temporal RSCs, one with correlation $\sigma = 0.3$ and the other with correlation $\sigma = 0.7$, and compute the fraction of infected nodes in the asymptotic state as a function of β_{\wedge} . As shown in Fig. 5.6,



Figure 5.6. Size of the infected population at the steady state as a function of β_{\triangle} for two different temporal correlations, $\sigma = 0.7$ and $\sigma = 0.3$. The critical value of the group infection rate β^c to enter the endemic state is lower for higher

The critical value of the group infection rate β_{\triangle}^c to enter the endemic state is lower for higher temporal correlation. Both curves display an abrupt transition as a function of β_{\triangle} . I set $\beta_{\parallel} = 0.85 \frac{\mu}{\langle k_{\parallel} \rangle}$ with $\langle k_{\parallel} \rangle = 12$ and $\langle k_{\triangle} \rangle = 5$.

in both cases the endemic phase is separated by an abrupt transition from the healthy region. The critical group infection rate for the transition to occur is higher in the first case.

I systematically investigate such phenomenon in Fig. 5.7, where I compute the critical group epidemic threshold as a function of σ . I observe that β^c_{Δ} decreases monotonically with the temporal correlation σ and it takes its minimum value for maximally correlated RSCs, corresponding to a static simplicial complex. Consistently with what was observed in Fig. 5.2 and Fig. 5.4, this suggests not only that group effects are weaker in temporal against a static setups, but that this is also the case the more diverse the temporal evolution of the system is.

I also note that the absence of a threshold β^c_{Δ} for values of temporal correlation below a critical σ^c , marked by a dashed vertical line, is due to the existence of a threshold of temporal correlation below which the transition to an endemic state is not possible, no matter the value of β_{Δ} .

5.5 Contagion on degree-heterogeneous temporal higher-order networks

In the previous section I investigated the effects of temporality in homogeneous simplicial complexes. I now turn my attention to the role of degree heterogeneity in temporal higher-order networks [161, 162, 186, 163].



Figure 5.7. Effect of temporal correlations in higher-order networks.

Critical group infection rate computed as a function of the temporal correlation σ . The critical value β_{Δ}^c is higher for decreasing values of σ , and the epidemic threshold disappears below a critical value of temporal correlation σ^c (grey line). This indicates that group effects are stronger in highly correlated higher-order networks. I set $\beta_{\parallel} = 0.85 \frac{\mu}{\langle k_{\parallel} \rangle}$, and each point in (c) was obtained by averaging over 100 RSCs with $\langle k_{\parallel} \rangle = 12$ and $\langle k_{\Delta} \rangle = 5$.

I generate scale-free (SF) simplicial complexes following a growth model introduced in [186], where both 1-simplices and 2-simplices follow a scale-free distribution, and where the sequences of k_{\parallel} and k_{\triangle} are maximally correlated. Next, I obtain a temporal sequence of SF simplicial complexes via recursively performing degree preserved rewiring at each time step such that the degree distribution of the simplices does not change. Desired values of temporal correlation can be achieved by suitably choosing the rewiring probability.

I simulate the epidemic process on top of two distinct sequences of SF simplicial complexes corresponding to the two extreme values of temporal correlation $\sigma_{max} = 1$ and $\sigma_{min} \approx 0$. For both configurations, I investigate two different scenarios of seeding infection, namely on the hub or on one of the leaves, and compute the fraction of the infected population in the long-time limit as a function of β_{\triangle} . As shown in Fig. 5.8a, for both hub and leaf cases, the critical value of β_{\triangle}^c to enter the endemic state is lower for higher values of temporal correlation, in agreement with what I found for homogeneous structures. Again, I only show the forward transition to the endemic state as the backward transition is not affected by temporality. As expected, seeding the infection on the hub enhances the epidemics. In particular, in the considered case, β_{\triangle}^c decreases by an order of magnitude when the infection is started on the best connected node of the network.

To properly quantify the effect of heterogeneity, I systematically compare the onset of the endemic state in the heterogeneous simplicial complex as a function



Figure 5.8. Effect of heterogeneity in higher-order networks.

(a): Fraction of the infected population on heterogeneous scale-free (SF) simplicial complexes (power-law exponents $\gamma_{\parallel} = 2.2$ and $\gamma_{\bigtriangleup} = 2.5$) as a function of β_{\bigtriangleup} for maximum (dashed lines) and minimum (solid lines) temporal correlation. I consider two different scenarios for initial infection: hub (red) and leaf (blue). High temporal correlation reduces the epidemic thresholds. (b): Epidemic thresholds as a function of β_{\parallel}^c and β_{\bigtriangleup}^c for heterogeneous and homogeneous simplicial complexes with the same number of interactions for the forward (solid curves) and backward (dashed curves) transition. On temporal SF complexes with no correlation, forward transition to the endemic state is possible for all considered values of β_{\parallel} upon increasing β_{\bigtriangleup} both in the hub (red) and leaf (blue) seeding scenarios, in contrast with RSCs (green) where the lack of temporal correlation prevents the onset of endemic phase entirely below a critical β_{\parallel} . For both SF and RSCs, the backward transition to infection-free state occurs upon decreasing β_{\bigtriangleup} , however a lower value of β_{\bigtriangleup} is required for SF complex as compared to RSCs. For panel (a), I set $\beta_{\parallel} = 0.25 \frac{\mu}{\langle k_{\circlearrowright} \rangle}$, for both (a) and (b), I set $\mu = 0.2$, $\langle k_{\parallel} \rangle = 10$, $\langle k_{\bigtriangleup} \rangle = 4$.

of both $\beta_{|}$ and β_{\triangle} against a homogeneous simplicial complex with the same number of 1- and 2-simplices. As shown in Fig. 5.8b, in uncorrelated temporal SF complex, for the forward transition, it is possible to reach the endemic state for all $\beta_{|}$ below a critical value upon increasing β_{\triangle} . This is in contrast with RSCs where, below a critical $\beta_{|}$, the lack of temporal correlation prevents the onset of the endemic phase entirely, as already observed in Fig. 5.2b. In such uncorrelated temporal case, for both SF and RSCs, the backward transition to infection-free state occurs upon decreasing β_{\triangle} , however a lower value of β_{\triangle} is required for SF complex as compared to RSCs. Homogeneous structures are the safer against contagion: when structural heterogeneity is present, starting the epidemic from a peripheral node will have a milder effect than if contagion begins from the hub, but the system is more prone to reach the endemic state compared to a homogeneous network with the same number of interactions.

5.6 Discussion

Motivated by data-driven investigations of careers of persistent teams in chapter 4, in this chapter I have introduced a modeling framework for spreading of ideas in temporally changing populations, which broadly speaking, lays the foundation for investigations in how ideas spread in science. In technical details, I investigated the effect of temporality and persistence of group interactions (or teams) by modeling the population as temporal higher-order networks. I focused primarily on the forward transition to the endemic state and showed that contagion processes behave remarkably differently on temporal and static finite-size homogeneous simplicial complexes. While in static networks the onset of the endemic state depends strongly on both β_{\parallel} and β_{\wedge} , in random temporal networks, where no correlations are present among timeconsecutive interaction structures, the effect of the higher-order contagion parameter is much weaker. This is linked to changes in the basins of attractions of the epidemic-free state, which shrinks fast for static structures when increasing the infectivity of the 2-simplices. As a consequence, temporality can have a direct impact on critical mass effects – already present in the static case [73]– by reshaping the basins of attractions of the system. In this scenario, a seed of infectious nodes of a fixed size can lead the system to both the endemic and epidemic-free states according to the temporal properties of its interactions. More in details, I investigated the effect of the initial infection size on the onset of the endemic state, finding that while for very small values of initial infection the onset of the epidemic is not impacted by group infectivity in both static and temporal simplicial complexes, a reasonable initial infection of size $\frac{1}{N}$ leads to striking differences between the two cases. Intermediate scenarios in the forward transition can be achieved on simplicial complexes with intermediate levels of temporal correlations. In contrast to the forward transition, I observed that the backward transition to infection-free state was unaffected by presence or absence of temporal correlations.

I also investigated the effect of degree heterogeneity on higher-order contagion. I confirmed that even in scale-free simplicial complexes, the absence of temporal correlations increases the infectivity required to achieve the endemic phase. However, in contrast to homogeneous simplicial complexes, in heterogeneous structures the lack of temporal correlations does not completely hinder the effect of group infectivity, and the endemic state can still be reached with a high enough value of β_{\wedge} . The parameter space associated to the endemic phase increases when the infection is seeded on a well-connected hub of the simplicial complex. However, even when the infection starts from a poorly connected node, the onset of the epidemics is always easier to achieve compared to an homogeneous simplicial complex with the same number of interactions.

As a possible limitation and future direction, I note that in this chapter, I did not attempt a calculation of the temporal correlation in scientific collaboration hypergraphs of different fields from OpenAlex data which I have used in Chapter 2 and 4. This could feed into the framework of this chapter allowing us to observe the effect of different discipline-specific temporal correlations in collaboration hypergraphs on the spreading of innovations. In the future, using my framework and comparing traditional fields like mathematics with fast changing fields like Machine learning and AI could also be interesting.

In the future, my temporal framework could be applied to investigate other dynamical processes recently extended beyond pairwise interactions, including opinion [170, 187], convention [171], and evolutionary dynamics [174]. Taken together, my analysis suggests the importance of considering temporality and persistence of interactions, a feature of many real-world systems, when investigating contagion processes on higher-order networks. As most higher-order social networks naturally evolve, with both pairwise and group interactions changing over time [182], my results suggest potential strategies to control contagion, by suitably tuning the temporal network structure. In context of the scientific ecosystem, this could enable achieving innovation faster by tuning the collaboration structure.

CHAPTER 5. MODEL OF SOCIAL CONTAGION IN A POPULATION WITH 86 EVOLVING GROUP STRUCTURE

CHAPTER 6

CONCLUSION

This thesis set out to quantify success in a data-driven fashion across a variety of domains, contributing to the emerging field known as Science of Success. It advances the understanding of success of single individuals by tracking them over larger career paths and extends this investigation to team success and their role in spreading innovations in later chapters. Unlike previous studies that focused on individual achievements or overlooked the significance of social networks, by looking at whole careers– for both individuals and teams– this research digs into the correlates of success in scientific and sporting contexts.

First, I investigated success in individual scientific careers from the lens of funding. I found that funding success depends not only on a scientist's work but also on where they stand in the wider network of science. Specifically, I discovered that early-career collaborations with U.S. institutions are crucial for securing major EU funding, revealing imbalances in EU-U.S. academic relationships, highlighting the role a scientists social network plays in funding success in the EU. This dependency suggests a potential risk of European research aligning too closely with American priorities, influenced more by strategic considerations than by scientific merit alone.

Secondly, the research extends to success patterns in individual careers in the sporting world, analyzing data from nearly a million chess players on Lichess.org. I found winning and losing streaks in chess, similar to those in scientific careers, indicating such dynamics may be common across multiple domains. By examining chess openings, I found players focus on fewer, more deeply understood strategies, developing a personal style over their career. Interestingly, players often do not end up using their most successful openings, especially experts, possibly due to the complex counterplays by skilled opponents. This study adds to the scientific understanding of human performance in chess.

Thirdly, I move beyond individual entities and investigated team careers in science, analyzing half a million teams. Teams with diverse disciplinary backgrounds tend to have shorter lifespans, highlighting the importance of achieving synergy early in long-term collaborations. High-impact periods for teams often cluster, similar to hot streaks in both chess and individual scientific careers. While multi-university teams have been seen as more successful since the 2000s, this is largely due to improved remote collaboration via ICT technologies. However, impactful teams balanced knowledge breadth well, contrasting with teams that either focused narrowly on topics or expanded knowledge too broadly, which often resulted in lower citation impacts. My research underscores the complexity of team dynamics in science, suggesting that composition may shape the success as well as longevity of collaborations. This work sets the stage for future research on how team composition, especially gender and ethnic diversity, impacts team success, and how strategic collaborations influence career trajectories in science.

Lastly, I explore how scientific ideas propagate within these temporally dynamic populations where members interact in groups. In particular, I explore how sustained, long-term interactions affect the spread of innovations using an agent-based model to social contagion. Using simulations on a temporally evolving population structures modelled via temporal simplicial-complexes, I found that group interactions bring about faster spreading of ideas but only when interactions are persistent. This study reveals the crucial role of persistent groups and temporal correlations in the contact structure in the spreading of innovations. By studying how ideas flow through groups that vary in their stability and member composition, we can uncover strategies to accelerate the diffusion process and potentially enhance the rate of collective scientific discovery.

In the future, methods I developed for quantifying human performance in chess could be used to study successful careers in other sports such as go, tennis, cricket etc. My results on success of scientific teams can be combined with prior theoretical work on team assembly to build a more complete framework of how teams form in science. Besides extracting persistent collaborations active over time, one can study team dynamics, how teams evolve and transform themselves and collectively adapt to a scientific ecosystem that constantly evolves in time. Furthermore, a cross-disciplinary analysis comparing teams in rapidly evolving fields like AI with those in more stable areas like mathematics, can guide tailored team formation strategies.

Future research could combine different themes of my thesis and examine the role of funding on team careers can highlight how external support mechanisms can promote meaningful persistent collaborations. Indeed, the scientific ecosystem is largely shaped in response to funding. From the perspective of individual careers, how persistent collaborations affect they way in which researchers navigate the knowledge landscape is also an intriguing question to address [134]. Long-term career outcomes of individuals who participate in high-impact, persistent collaborative experiences can be compared to those who work in more transient or less successful teams. Notably, early-career collaborations with elite scientists predict subsequent career success [135]. Investigating this aspect further could provide insights into how early career researchers can strategically navigate collaborations and switch between teams to enhance their professional growth and impact. Lastly, a data-driven investigation into the roles that teams play in spreading ideas within scientific ecosystems could be used to validate the theoretical framework developed in the final chapter of this thesis. This approach could also help identify the types of teams that are most effective at maximizing the spread of innovative ideas.

In summary, this thesis contributes to the emerging field of Science of Success by identifying key factors that drive successful careers in both individual and team settings, in science and sports. It highlights the critical role of social networks, collaboration dynamics, and long-term interactions in shaping career trajectories and success. By integrating these insights, this work offers valuable guidance for funding bodies, institutions, and individuals aiming to foster sustained success in competitive intellectual fields. Overall, this thesis focuses on success in time, characterising the careers of successful individuals and teams and quantifies the difference in behaviour of successful and the less successful scientists and chess players.

I hope this work opens doors for future research that employs the detailed data analysis presented in this thesis to understand the factors leading to successful careers in intellectual domains. In Science, it can inform various stake holders such as funding agencies and institutions on which individual careers and teams to support. In Sports, it can serve as a roadmap for developing players on how to be shape their careers to set themselves up for success.

BIBLIOGRAPHY

- [1] Michael Lewis. *Moneyball: The art of winning an unfair game*. WW Norton & Company, 2004.
- [2] Pierre Bourdieu. *The field of cultural production: Essays on art and literature*. Columbia University Press, 1993.
- [3] Brian Uzzi. A social network's changing statistical properties and the quality of human innovation. *Journal of Physics A: Mathematical and Theoretical*, 41(22):224023, 2008.
- [4] Qing Ke, Emilio Ferrara, Filippo Radicchi, and Alessandro Flammini. Defining and identifying sleeping beauties in science. *Proceedings of the National Academy of Sciences*, 112(24):7426–7431, 2015.
- [5] Albert-László Barabási. *The formula: The universal laws of success*. Hachette UK, 2018.
- [6] Samuel P Fraiberger, Roberta Sinatra, Magnus Resch, Christoph Riedl, and Albert-László Barabási. Quantifying reputation and success in art. *Science*, 362(6416):825–829, 2018.
- [7] Maximilian Schich and Isabel Meirelles. Arts, humanities and complex networks: Introduction. *Leonardo*, 49(5):445–445, 2016.
- [8] Burcu Yucesoy and Albert-László Barabási. Untangling performance from success. *EPJ Data Science*, 5(1):1–10, 2016.
- [9] Daniel J Boorstin. *The image: A guide to pseudo-events in America*. Vintage, 1992.
- [10] Amy Argetsinger. Famesque: Amy argetsinger on celebrities famous for being famous. *Washington Post*, 2013.

- [11] Neal Gabler. Toward a new definition of celebrity. USC Annenberg: The Norman Lear Center, 2001.
- [12] Arnout Van de Rijt, Eran Shor, Charles Ward, and Steven Skiena. Only 15 minutes? the social stratification of fame in printed media. *American Sociological Review*, 78(2):266–289, 2013.
- [13] Openalex: A fully-open index of scholarly works, authors, venues, institutions, and concepts. https://openalex.org/.
- [14] Dalmeet Singh Chawla. Massive open index of scholarly papers launches. *Nature*, 2022.
- [15] Diana Hicks, Paul Wouters, Ludo Waltman, Sarah De Rijcke, and Ismael Rafols. Bibliometrics: the leiden manifesto for research metrics. *Nature News*, 520(7548):429, 2015.
- [16] An Zeng, Zhesi Shen, Jianlin Zhou, Jinshan Wu, Ying Fan, Yougui Wang, and H Eugene Stanley. The science of science: From the perspective of complex systems. *Physics reports*, 714:1–73, 2017.
- [17] Aaron Clauset, Daniel B Larremore, and Roberta Sinatra. Data-driven predictions in the science of science. *Science*, 355(6324):477–480, 2017.
- [18] Cassidy R Sugimoto and Vincent Larivière. *Measuring research: What everyone needs to know*. Oxford University Press, 2018.
- [19] Santo Fortunato, Carl T Bergstrom, Katy Börner, James A Evans, Dirk Helbing, Staša Milojević, Alexander M Petersen, Filippo Radicchi, Roberta Sinatra, Brian Uzzi, et al. Science of science. *Science*, 359(6379):eaao0185, 2018.
- [20] Lu Liu, Benjamin F Jones, Brian Uzzi, and Dashun Wang. Data, measurement and empirical methods in the science of science. *Nature human behaviour*, 7(7):1046–1058, 2023.
- [21] Roberta Sinatra, Dashun Wang, Pierre Deville, Chaoming Song, and Albert-László Barabási. Quantifying the evolution of individual scientific impact. *Science*, 354(6312), 2016.
- [22] Lu Liu, Yang Wang, Roberta Sinatra, C Lee Giles, Chaoming Song, and Dashun Wang. Hot streaks in artistic, cultural, and scientific careers. *Nature*, 559(7714):396–399, 2018.

- [23] Lu Liu, Nima Dehmamy, Jillian Chown, C Lee Giles, and Dashun Wang. Understanding the onset of hot streaks across artistic, cultural, and scientific careers. *Nature communications*, 12(1):5392, 2021.
- [24] Vincent Larivière, Chaoqun Ni, Yves Gingras, Blaise Cronin, and Cassidy R Sugimoto. Bibliometrics: Global gender disparities in science. *Nature*, 504(7479):211–213, 2013.
- [25] Junming Huang, Alexander J Gates, Roberta Sinatra, and Albert-László Barabási. Historical comparison of gender inequality in scientific careers across countries and disciplines. *Proceedings of the National Academy of Sciences*, 117(9):4609–4616, 2020.
- [26] Maxwell A Bertolero, Jordan D Dworkin, Sophia U David, Claudia López Lloreda, Pragya Srivastava, Jennifer Stiso, Dale Zhou, Kafui Dzirasa, Damien A Fair, Antonia N Kaczkurkin, et al. Racial and ethnic imbalance in neuroscience reference lists and intersections with gender. *BioRxiv*, pages 2020–10, 2020.
- [27] Dennis L. Murray, Douglas Morris, Claude Lavoie, Peter R. Leavitt, Hugh MacIsaac, Michael E. J. Masson, and Marc-Andre Villard. Bias in research grant evaluation has dire consequences for small universities. *PloS One*, 11(6):1–19, 06 2016.
- [28] Sam Zhang, K Hunter Wapman, Daniel B Larremore, and Aaron Clauset. Labor advantages drive the greater productivity of faculty at elite universities. *Science Advances*, 8(46):eabq7056, 2022.
- [29] Alexander Krauss, Lluís Danús, and Marta Sales-Pardo. Early-career factors largely determine the future impact of prominent researchers: evidence across eight scientific fields. *Scientific Reports*, 13(1):18794, 2023.
- [30] Pierre Deville, Dashun Wang, Roberta Sinatra, Chaoming Song, Vincent D Blondel, and Albert-László Barabási. Career on the move: Geography, stratification and scientific impact. *Scientific reports*, 4(1):1–7, 2014.
- [31] An Zeng, Zhesi Shen, Jianlin Zhou, Ying Fan, Zengru Di, Yougui Wang, H Eugene Stanley, and Shlomo Havlin. Increasing trend of scientists to switch between topics. *Nature communications*, 10(1):1–11, 2019.
- [32] Staša Milojević, Filippo Radicchi, and John P Walsh. Changing demographics of scientific careers: The rise of the temporary workforce. *Proceedings of the National Academy of Sciences*, 115(50):12616–12623, 2018.

- [33] T. K. Woodruff Y. Ma, D. F. M. Oliveira and B. Uzzi. Women who win prizes get less money and prestige. *Nature*, 565:287–288, 2019.
- [34] J Mervis. In effort to understand continuing racial disparities, NIH to test for bias in study sections. *Science*, 2016.
- [35] Lindell Bromham, Russell Dinnage, and Xia Hua. Interdisciplinary research has consistently lower funding success. *Nature*, 534(7609):684–687, Jun 2016.
- [36] Thijs Bol, Mathijs de Vaan, and Arnout van de Rijt. The matthew effect in science funding. *Proceedings of the National Academy of Sciences*, 115(19):4887–4890, 2018.
- [37] Athen Ma, Raúl J. Mondragón, and Vito Latora. Anatomy of funded research in science. *Proceedings of the National Academy of Sciences*, 112(48):14760–14765, 2015.
- [38] Michael Szell and Roberta Sinatra. Research funding goes to rich clubs. *Proceedings of the National Academy of Sciences*, 112(48):14749–14750, 2015.
- [39] K. Hunter Wapman, Sam Zhang, Aaron Clauset, and Daniel B. Larremore. Quantifying hierarchy and dynamics in US faculty hiring and retention. *Nature*, 610(7930):120–127, Oct 2022.
- [40] Alan Michael Nevill, Greg Atkinson, and Mike Hughes. Twenty-five years of sport performance research in the journal of sports sciences. *Journal of Sports Sciences*, 26:413 – 426, 2008.
- [41] Damon PS Andrew, Paul M Pedersen, and Chad D McEvoy. *Research methods and design in sport management*. Human Kinetics, 2019.
- [42] Filippo Radicchi. Who is the best player ever? a complex network analysis of the history of professional tennis. *PloS one*, 6(2):e17249, 2011.
- [43] Merim Bilalić, Peter McLeod, and Fernand Gobet. Does chess need intelligence?—a study with young chess players. *Intelligence*, 35(5):457–470, 2007.
- [44] Neil Charness, Michael Tuffiash, Ralf Krampe, Eyal Reingold, and Ekaterina Vasyukova. The role of deliberate practice in chess expertise. *Applied Cognitive Psychology*, 19(2):151–165, 2005.
- [45] Guillermo Campitelli and Fernand Gobet. Deliberate practice: Necessary but not sufficient. *Current directions in psychological science*, 20(5):280–285, 2011.
- [46] David Z Hambrick, Frederick L Oswald, Erik M Altmann, Elizabeth J Meinz, Fernand Gobet, and Guillermo Campitelli. Deliberate practice: Is that all it takes to become an expert? *Intelligence*, 45:34–45, 2014.
- [47] Roger Guimera, Brian Uzzi, Jarrett Spiro, and Luis A Nunes Amaral. Team assembly mechanisms determine collaboration network structure and team performance. *Science*, 308(5722):697–702, 2005.
- [48] Stefan Wuchty, Benjamin F Jones, and Brian Uzzi. The increasing dominance of teams in production of knowledge. *Science*, 316(5827):1036–1039, 2007.
- [49] Laura Hunter and Erin Leahey. Collaborative research in sociology: Trends and contributing factors. *The American Sociologist*, 39(4):290–306, 2008.
- [50] Benjamin F Jones, Stefan Wuchty, and Brian Uzzi. Multi-university research teams: Shifting impact, geography, and stratification in science. *science*, 322(5905):1259–1262, 2008.
- [51] Stephanie Teasley and Steven Wolinsky. Scientific collaborations at a distance, 2001.
- [52] Brian Uzzi, Satyam Mukherjee, Michael Stringer, and Ben Jones. Atypical combinations and scientific impact. *Science*, 342(6157):468–472, 2013.
- [53] Alina Lungeanu, Yun Huang, and Noshir S Contractor. Understanding the assembly of interdisciplinary teams and its impact on performance. *Journal of informetrics*, 8(1):59–70, 2014.
- [54] Satyam Mukherjee, Yun Huang, Julia Neidhardt, Brian Uzzi, and Noshir Contractor. Prior shared success predicts victory in team competitions. *Nature human behaviour*, 3(1):74–81, 2019.
- [55] Kara L Hall, Amanda L Vogel, Grace C Huang, Katrina J Serrano, Elise L Rice, Sophia P Tsakraklides, and Stephen M Fiore. The science of team science: A review of the empirical evidence and research gaps on collaboration in science. *American psychologist*, 73(4):532, 2018.

- [56] Lingfei Wu, Dashun Wang, and James A Evans. Large teams develop and small teams disrupt science and technology. *Nature*, 566(7744):378–382, 2019.
- [57] David A Harrison, Susan Mohammed, Joseph E McGrath, Anna T Florey, and Scott W Vanderstoep. Time matters in team performance: Effects of member familiarity, entrainment, and task discontinuity on speed and quality. *Personnel Psychology*, 56(3):633–669, 2003.
- [58] Bedoor K AlShebli, Talal Rahwan, and Wei Lee Woon. The preeminence of ethnic diversity in scientific collaboration. *Nature communications*, 9(1):1– 10, 2018.
- [59] Staša Milojević. Principles of scientific research team formation and evolution. *Proceedings of the National Academy of Sciences*, 111(11):3984–3989, 2014.
- [60] Itai Yanai and Martin J Lercher. It takes two to think. *Nature Biotechnology*, pages 1–2, 2024.
- [61] Anne Q Hoy. Science diplomacy leverages alliances to build global bridges. *Science*, 365(6456):875–876, 2019.
- [62] Alexander Michael Petersen. Quantifying the impact of weak, strong, and super ties in scientific careers. *Proceedings of the National Academy of Sciences*, 112(34):E4671–E4680, 2015.
- [63] Junwan Liu, Xiaofei Guo, Shuo Xu, Yi Bu, Cassidy R Sugimoto, Vincent Larivière, Yinglu Song, and Honghao Zhou. Understanding superpartnerships in scientific collaboration: Evidence from the field of economics. *Journal of the Association for Information Science and Technology*, 2024.
- [64] Gangmin Son, Jinhyuk Yun, and Hawoong Jeong. Untangling pair synergy in the evolution of collaborative scientific impact. *EPJ Data Science*, 12(1):62, 2023.
- [65] Romualdo Pastor-Satorras, Claudio Castellano, Piet Van Mieghem, and Alessandro Vespignani. Epidemic processes in complex networks. *Re*views of modern physics, 87(3):925, 2015.
- [66] Douglas Guilbeault, Joshua Becker, and Damon Centola. Complex contagions: A decade in review. *Complex spreading phenomena in social systems: Influence and contagion in real-world social networks*, pages 3–25, 2018.

- [67] Damon Centola and Michael Macy. Complex contagions and the weakness of long ties. *American journal of Sociology*, 113(3):702–734, 2007.
- [68] Damon Centola. The spread of behavior in an online social network experiment. *science*, 329(5996):1194–1197, 2010.
- [69] Johan Ugander, Lars Backstrom, Cameron Marlow, and Jon Kleinberg. Structural diversity in social contagion. *Proceedings of the national academy* of sciences, 109(16):5962–5966, 2012.
- [70] Lilian Weng, Alessandro Flammini, Alessandro Vespignani, and Fillipo Menczer. Competition among memes in a world with limited attention. *Scientific reports*, 2(1):335, 2012.
- [71] Márton Karsai, Gerardo Iniguez, Kimmo Kaski, and János Kertész. Complex contagion process in spreading of online innovation. *Journal of The Royal Society Interface*, 11(101):20140694, 2014.
- [72] Bjarke Mønsted, Piotr Sapieżyński, Emilio Ferrara, and Sune Lehmann. Evidence of complex contagion of information in social media: An experiment using twitter bots. *PloS one*, 12(9):e0184148, 2017.
- [73] Iacopo Iacopini, Giovanni Petri, Alain Barrat, and Vito Latora. Simplicial models of social contagion. *Nat. Commun.*, 10(1):2485, 2019.
- [74] Giorgio Parisi. Balance research funds across europe. *Nature*, 530(7588):33–33, Feb 2016.
- [75] Manlio De Domenico and Alex Arenas. Eu cash goes to the sticky and attractive. *Nature*, 531(7596):580–580, Mar 2016.
- [76] EU grants: List of participating countries. https://ec.europa.eu/info/ funding-tenders/opportunities/docs/2021-2027/common/guidance/ list-3rd-country-participation_horizon-euratom_en.pdf.
- [77] D. S. Chawla. Massive open index of scholarly papers launches. *Nature*, Jan 2022.
- [78] Panel structure for ERC calls 2021 and 2022. https://erc.europa. eu/sites/default/files/document/file/ERC_Panel_structure_2021_ 2022.pdf.
- [79] Sandeep Chowdhary, Nicolò Defenu, Federico Musciotto, and Federico Battiston. Funding bias: nurture european researchers' independence. *Nature*, 616(7955):33, 2023.

- [80] C. A. M. L. Porta and S. Zapperi. America's top universities reap the benefit of italian-trained scientists. *Nature Italy*, 2022.
- [81] ERC scientific council decides changes to the evaluation forms and processes for the 2024 calls. https://erc.europa.eu/newsevents/news/ erc-scientific-council-decides-changes-evaluation-forms-and-processes-2024-call html, 2022.
- [82] Ana L Schaigorodsky, Juan I Perotti, and Orlando V Billoni. Memory and long-range correlations in chess games. *Physica A: Statistical Mechanics and its Applications*, 394:304–311, 2014.
- [83] Bernd Blasius and Ralf Tönjes. Zipf's law in the popularity distribution of chess openings. *Phys. Rev. Lett.*, 103(21):218701, 2009.
- [84] Ramiz Arabacı. An investigation into the openings used by top 100 chess players. *International Journal of Performance Analysis in Sport*, 6(1):149–160, 2006.
- [85] Claude Elwood Shannon. A mathematical theory of communication. *ACM SIGMOBILE mobile computing and communications review*, 5(1):3–55, 2001.
- [86] Mark E Glickman. Example of the glicko-2 system. *Boston University*, pages 1–6, 2012.
- [87] Thomas Gilovich, Robert Vallone, and Amos Tversky. The hot hand in basketball: On the misperception of random sequences. *Cognitive Psychology*, 17(3):295–314, 1985.
- [88] Joshua B Miller and Adam Sanjurjo. Surprised by the hot hand fallacy? a truth in the law of small numbers. *Econometrica*, 86(6):2019–2047, 2018.
- [89] Janet A Young and Michelle D Pain. The zone: Evidence of a universal phenomenon for athletes across sports. *Athletic Insight: the online journal of sport psychology*, 1(3):21–30, 1999.
- [90] Michael Murphy and Rhea A White. In the zone: Transcendent experience in sports. 2011.
- [91] Steven D Hales. An epistemologist looks at the hot hand in sports. *Journal* of the Philosophy of Sport, 26(1):79–87, 1999.

- [92] Michael Bar-Eli, Simcha Avugos, and Markus Raab. Twenty years of "hot hand" research: Review and critique. *Psychology of Sport and Exercise*, 7(6):525–553, 2006.
- [93] T.S. Kuhn. The essential tension: Selected studies in scientific tradition and change. 2011.
- [94] James G March. Exploration and exploitation in organizational learning. *Organization Science*, 2(1):71–87, 1991.
- [95] Luca Pappalardo, Filippo Simini, Salvatore Rinzivillo, Dino Pedreschi, Fosca Giannotti, and Albert-László Barabási. Returners and explorers dichotomy in human mobility. *Nature Communications*, 6(1):1–8, 2015.
- [96] Iacopo Iacopini, Staša Milojević, and Vito Latora. Network dynamics of innovation processes. *Physical Review Letters*, 120(4):048301, 2018.
- [97] Adrianus Dingeman de Groot. Thought and choice in chess. 1978.
- [98] Herbert Simon and William Chase. Skill in chess. pages 175–188, 1988.
- [99] William G. Chase and Herbert A. Simon. Perception in chess. Cogn. Psychol., 4:55–81, 1973.
- [100] Herbert A. Simon and Kevin Michael Gilmartin. A simulation of memory for chess positions. *Cognitive Psychology*, 5:29–46, 1973.
- [101] Milán Janosov, Federico Battiston, and Roberta Sinatra. Success and luck in creative careers. *EPJ Data Science*, 9:1–12, 2020.
- [102] Giordano De Marzo and Vito DP Servedio. Quantifying the complexity and similarity of chess openings using online chess community data. *arXiv preprint arXiv:2206.14312, 2022.*
- [103] Allen Newell, Herbert Alexander Simon, et al. Human problem solving. 104(9), 1972.
- [104] Peter M Todd and Gerd Gigerenzer. Précis of simple heuristics that make us smart. *Behavioral and brain sciences*, 23(5):727–741, 2000.
- [105] Gerd Gigerenzer and Wolfgang Gaissmaier. Heuristic decision making. *Annual review of psychology*, 62(1):451–482, 2011.
- [106] Zoltán Medvegy, Markus Raab, Kata Tóth, Gergely Csurilla, and Tamás Sterbenz. When do expert decision makers trust their intuition? *Applied Cognitive Psychology*, 2022.

- [107] Joseph G Johnson and Markus Raab. Take the first: Option-generation and resulting choices. *Organizational behavior and human decision processes*, 91(2):215–229, 2003.
- [108] Tao Jia, Dashun Wang, and Boleslaw K Szymanski. Quantifying patterns of research-interest evolution. *Nature Human Behaviour*, 1(4):1–7, 2017.
- [109] David Hsiehchen, Magdalena Espinoza, and Antony Hsieh. Multinational teams and diseconomies of scale in collaborative research. *Science advances*, 1(8):e1500211, 2015.
- [110] Mario Coccia and Lili Wang. Evolution and convergence of the patterns of international scientific collaboration. *Proceedings of the National Academy* of Sciences, 113(8):2057–2061, 2016.
- [111] Ali Gazni, Cassidy R Sugimoto, and Fereshteh Didegah. Mapping world scientific collaboration: Authors, institutions, and countries. *Journal of the American Society for Information Science and Technology*, 63(2):323–335, 2012.
- [112] Mathias Wullum Nielsen, Sharla Alegria, Love Börjeson, Henry Etzkowitz, Holly J Falk-Krzesinski, Aparna Joshi, Erin Leahey, Laurel Smith-Doerr, Anita Williams Woolley, and Londa Schiebinger. Gender diversity leads to better science. *Proceedings of the National Academy of Sciences*, 114(8):1740–1742, 2017.
- [113] An Zeng, Ying Fan, Zengru Di, Yougui Wang, and Shlomo Havlin. Fresh teams are associated with original and multidisciplinary research. *Nature Human Behaviour*, pages 1–9, 2021.
- [114] Federico Musciotto, Federico Battiston, and Rosario N Mantegna. Detecting informative higher-order interactions in statistically validated hypergraphs. *arXiv preprint arXiv:2103.16484*, 2021.
- [115] Openalex: Author disambiguation. https://docs.openalex.org/ api-entities/authors/author-disambiguation.
- [116] Federico Musciotto, Federico Battiston, and Rosario N Mantegna. Identifying maximal sets of significantly interacting nodes in higher-order networks. arXiv preprint arXiv:2209.12712, 2022.
- [117] Federico Battiston, Giulia Cencetti, Iacopo Iacopini, Vito Latora, Maxime Lucas, Alice Patania, Jean-Gabriel Young, and Giovanni Petri. Networks beyond pairwise interactions: structure and dynamics. *Phys. Rep.*, 874:1– 92, 2020.

- [118] Minghui Wang, Yongzhong Zhao, and Bin Zhang. Efficient test and visualization of multi-set intersections. *Scientific reports*, 5(1):16923, 2015.
- [119] Openalex: Concepts. https://docs.openalex.org/api-entities/ concepts.
- [120] Milán Janosov, Federico Battiston, and Roberta Sinatra. Success and luck in creative careers. *EPJ Data Science*, 9(1):9, 2020.
- [121] Filippo Radicchi, Santo Fortunato, and Claudio Castellano. Universality of citation distributions: Toward an objective measure of scientific impact. *Proceedings of the National Academy of Sciences*, 105(45):17268–17272, 2008.
- [122] Oliver E Williams, Lucas Lacasa, and Vito Latora. Quantifying and predicting success in show business. *Nature communications*, 10(1):1–8, 2019.
- [123] Alan Porter and Ismael Rafols. Is science becoming more interdisciplinary? measuring and mapping six research fields over time. *Scientometrics*, 81(3):719–745, 2009.
- [124] Richard Van Noorden et al. Interdisciplinary research by the numbers. *Nature*, 525(7569):306–307, 2015.
- [125] Longqi Yang, David Holtz, Sonia Jaffe, Siddharth Suri, Shilpi Sinha, Jeffrey Weston, Connor Joyce, Neha Shah, Kevin Sherman, Brent Hecht, et al. The effects of remote work on collaboration among information workers. *Nature human behaviour*, 6(1):43–54, 2022.
- [126] Yiling Lin, Carl Benedikt Frey, and Lingfei Wu. Remote collaboration fuses fewer breakthrough ideas. *Nature*, 623(7989):987–991, 2023.
- [127] Yang Yang, Tanya Y Tian, Teresa K Woodruff, Benjamin F Jones, and Brian Uzzi. Gender-diverse teams produce more novel and higherimpact scientific ideas. *Proceedings of the National Academy of Sciences*, 119(36):e2200841119, 2022.
- [128] Diego Kozlowski, Dakota S Murray, Alexis Bell, Will Hulsey, Vincent Larivière, Thema Monroe-White, and Cassidy R Sugimoto. Avoiding bias when inferring race using name-based approaches. *Plos one*, 17(3):e0264270, 2022.
- [129] Jeffrey W Lockhart, Molly M King, and Christin Munsch. Name-based demographic inference and the unequal distribution of misrecognition. *Nature Human Behaviour*, 7(7):1084–1095, 2023.

- [130] Staša Milojević. Modes of collaboration in modern science: Beyond power laws and preferential attachment. *Journal of the american society for Information science and technology*, 61(7):1410–1423, 2010.
- [131] Mirta Galesic, Daniel Barkoczi, Andrew M Berdahl, Dora Biro, Giuseppe Carbone, Ilaria Giannoccaro, Robert L Goldstone, Cleotilde Gonzalez, Anne Kandler, Albert B Kao, et al. Beyond collective intelligence: Collective adaptation. *Journal of the Royal Society interface*, 20(200):20220736, 2023.
- [132] Lindell Bromham, Russell Dinnage, and Xia Hua. Interdisciplinary research has consistently lower funding success. *Nature*, 534(7609):684–687, 2016.
- [133] Sandeep Chowdhary, Nicolò Defenu, Federico Musciotto, and Federico Battiston. Dependency of ERC-funded research on US collaborations. *Nature Physics*, 19(12):1746–1749, 2023.
- [134] Chakresh Kumar Singh, Liubov Tupikina, Fabrice Lécuyer, Michele Starnini, and Marc Santolini. Charting mobility patterns in the scientific knowledge landscape. *EPJ Data Science*, 13(1):12, 2024.
- [135] Weihua Li, Tomaso Aste, Fabio Caccioli, and Giacomo Livan. Early coauthorship with top scientists predicts success in academic careers. *Nature communications*, 10(1):5170, 2019.
- [136] Matt J Keeling and Pejman Rohani. *Modeling infectious diseases in humans and animals*. Princeton university press, 2011.
- [137] William Goffman and Vaun A Newill. Generalization of epidemic theory: An application to the transmission of ideas. *Nature*, 204(4955):225–228, 1964.
- [138] Daryl J Daley and David G Kendall. Epidemics and rumours. 1964.
- [139] Stefano Boccaletti, Vito Latora, Yamir Moreno, Martin Chavez, and D-U Hwang. Complex networks: Structure and dynamics. *Phys. Rep.*, 424(4-5):175–308, 2006.
- [140] Alain Barrat, Marc Barthelemy, and Alessandro Vespignani. *Dynamical processes on complex networks*. Cambridge university press, 2008.
- [141] Mason A Porter and James P Gleeson. Dynamical systems on networks. *Front. Appl. Dyn. Syst.: Rev. Tutor.*, 4, 2016.

- [142] Guilherme Ferraz de Arruda, Francisco A Rodrigues, and Yamir Moreno. Fundamentals of spreading processes in single and multilayer complex networks. *Phys. Rep.*, 756:1–59, 2018.
- [143] Romualdo Pastor-Satorras and Alessandro Vespignani. Epidemic spreading in scale-free networks. *Phys. Rev. Lett.*, 86(14):3200, 2001.
- [144] Duncan J Watts and Peter Sheridan Dodds. Influentials, networks, and public opinion formation. *J. Consum. Res.*, 34(4):441–458, 2007.
- [145] Sune Lehmann and Yong-Yeol Ahn. *Complex spreading phenomena in social systems*. Springer, 2018.
- [146] Duncan J Watts. A simple model of global cascades on random networks. In *The Structure and Dynamics of Networks*, pages 497–502. Princeton University Press, 2011.
- [147] Réka Albert and Albert-László Barabási. Statistical mechanics of complex networks. *Rev. Mod. Phys.*, 74(1):47, 2002.
- [148] V. Latora, V. Nicosia, and G. Russo. *Complex Networks: Principles, Methods and Applications*. Cambridge University Press, 2017.
- [149] Renaud Lambiotte, Martin Rosvall, and Ingo Scholtes. From networks to optimal higher-order models of complex systems. *Nat. Phys.*, 15(4):313– 320, 2019.
- [150] Leo Torres, Ann S Blevins, Danielle S Bassett, and Tina Eliassi-Rad. The why, how, and when of representations for complex systems. *arXiv preprint arXiv*:2006.02870, 2020.
- [151] Christian Bick, Elizabeth Gross, Heather A Harrington, and Michael T Schaub. What are higher-order networks? *arXiv preprint arXiv:2104.11329*, 2021.
- [152] Christian Kuehn and Christian Bick. A universal route to explosive phenomena. *Sci. Adv.*, 7(16):eabe3824, 2021.
- [153] Li Chen, Fakhteh Ghanbarnejad, and Dirk Brockmann. Fundamental properties of cooperative contagion processes. *New J. Phys.*, 19(10):103041, 2017.
- [154] J Gómez-Gardeñes, L Lotero, SN Taraskin, and FJ Pérez-Reche. Explosive contagion in networks. *Sci. Rep.*, 6(1):1–9, 2016.

- [155] David Soriano-Paños, Quantong Guo, Vito Latora, and Jesús Gómez-Gardeñes. Explosive transitions induced by interdependent contagionconsensus dynamics in multiplex networks. *Phys. Rev. E*, 99(6):062311, 2019.
- [156] David Soriano-Paños, Fakhteh Ghanbarnejad, Sandro Meloni, and Jesús Gómez-Gardeñes. Markovian approach to tackle the interaction of simultaneous diseases. *Phys. Rev. E*, 100(6):062308, 2019.
- [157] Alain Barrat, Guilherme Ferraz de Arruda, Iacopo Iacopini, and Yamir Moreno. Social contagion on higher-order structures. *arXiv preprint arXiv:*2103.03709, 2021.
- [158] Joan T Matamalas, Sergio Gómez, and Alex Arenas. Abrupt phase transition of epidemic spreading in simplicial complexes. *Phys. Rev. Res.*, 2(1):012049, 2020.
- [159] Giulio Burgio, Alex Arenas, Sergio Gómez, and Joan T Matamalas. Network clique cover approximation to analyze complex contagions through group interactions. *Commun. Phys.*, 4, 111, 2021.
- [160] Guillaume St-Onge, Iacopo Iacopini, Vito Latora, Alain Barrat, Giovanni Petri, Antoine Allard, and Laurent Hébert-Dufresne. Influential groups for seeding and sustaining hypergraph contagions. *arXiv preprint arXiv*:2105.07092, 2021.
- [161] Guilherme Ferraz de Arruda, Giovanni Petri, and Yamir Moreno. Social contagion models on hypergraphs. *Phys. Rev. Res.*, 2(2):023032, 2020.
- [162] Nicholas W Landry and Juan G Restrepo. The effect of heterogeneity on hypergraph contagion models. *Chaos: Interdiscip. J. Nonlinear Sci.*, 30(10):103117, 2020.
- [163] Guilherme Ferraz de Arruda, Michele Tizzani, and Yamir Moreno. Phase transitions and stability of dynamical processes on hypergraphs. *Commun. Phys.*, 4(1):1–9, 2021.
- [164] Christian Bick, Peter Ashwin, and Ana Rodrigues. Chaos in generically coupled phase oscillator networks with nonpairwise interactions. *Chaos: Interdiscip. J. Nonlinear Sci.*, 26(9):094814, 2016.
- [165] Per Sebastian Skardal and Alex Arenas. Abrupt desynchronization and extensive multistability in globally coupled oscillator simplexes. *Phys. Rev. Lett.*, 122(24):248301, 2019.

- [166] Ana P Millán, Joaquín J Torres, and Ginestra Bianconi. Explosive higherorder kuramoto dynamics on simplicial complexes. *Physical Review Letters*, 124(21):218301, 2020.
- [167] Maxime Lucas, Giulia Cencetti, and Federico Battiston. Multiorder laplacian for synchronization in higher-order networks. *Physical Review Research*, 2(3):033410, 2020.
- [168] Timoteo Carletti, Federico Battiston, Giulia Cencetti, and Duccio Fanelli. Random walks on hypergraphs. *Phys. Rev. E*, 101(2):022308, 2020.
- [169] Michael T Schaub, Austin R Benson, Paul Horn, Gabor Lippner, and Ali Jadbabaie. Random walks on simplicial complexes and the normalized hodge 1-laplacian. SIAM Review, 62(2):353–391, 2020.
- [170] Leonie Neuhäuser, Andrew Mellor, and Renaud Lambiotte. Multibody interactions and nonlinear consensus dynamics on networked systems. *Phys. Rev. E*, 101(3):032310, 2020.
- [171] Iacopo Iacopini, Giovanni Petri, Andrea Baronchelli, and Alain Barrat. Vanishing size of critical mass for tipping points in social convention. *arXiv preprint arXiv:2103.10411*, 2021.
- [172] Eyal Bairey, Eric D Kelsic, and Roy Kishony. High-order species interactions shape ecosystem diversity. *Nat. Commun.*, 7(1):1–7, 2016.
- [173] Jacopo Grilli, György Barabás, Matthew J Michalska-Smith, and Stefano Allesina. Higher-order interactions stabilize dynamics in competitive network models. *Nature*, 548(7666):210–213, 2017.
- [174] Unai Alvarez-Rodriguez, Federico Battiston, Guilherme Ferraz de Arruda, Yamir Moreno, Matjaž Perc, and Vito Latora. Evolutionary dynamics of higher-order interactions in social networks. *Nature Human Behaviour*, 5(5):586–595, 2021.
- [175] Luis EC Rocha, Fredrik Liljeros, and Petter Holme. Simulated epidemics in an empirical spatiotemporal network of 50,185 sexual contacts. *PLoS Comput. Biol.*, 7(3):e1001109, 2011.
- [176] Márton Karsai, Mikko Kivelä, Raj Kumar Pan, Kimmo Kaski, János Kertész, A-L Barabási, and Jari Saramäki. Small but slow world: How network topology and burstiness slow down spreading. *Phys. Rev. E*, 83(2):025102, 2011.

- [177] Petter Holme and Jari Saramäki. Temporal networks. *Physics reports*, 519(3):97–125, 2012.
- [178] Nicola Perra, Bruno Gonçalves, Romualdo Pastor-Satorras, and Alessandro Vespignani. Activity driven modeling of time varying networks. *Sci. Rep.*, 2(1):1–7, 2012.
- [179] Eugenio Valdano, Luca Ferreri, Chiara Poletto, and Vittoria Colizza. Analytical computation of the epidemic threshold on temporal networks. *Phys. Rev. X*, 5(2):021005, 2015.
- [180] Naoki Masuda and Petter Holme. *Introduction to temporal network epidemi*ology. Springer, 2017.
- [181] Andreas Koher, Hartmut HK Lentz, James P Gleeson, and Philipp Hövel. Contact-based model for epidemic spreading on temporal networks. *Phys. Rev. X*, 9(3):031017, 2019.
- [182] Giulia Cencetti, Federico Battiston, Bruno Lepri, and Márton Karsai. Temporal properties of higher-order interactions in social networks. *Scientific reports*, 11(1):1–10, 2021.
- [183] Guillaume St-Onge, Hanlin Sun, Antoine Allard, Laurent Hébert-Dufresne, and Ginestra Bianconi. Bursty exposure on higherorder networks leads to nonlinear infection kernels. arXiv preprint arXiv:2101.07229, 2021.
- [184] Sergio Gómez, Alexandre Arenas, Javier Borge-Holthoefer, Sandro Meloni, and Yamir Moreno. Discrete-time markov chain approach to contactbased disease spreading in complex networks. *Europhys. Lett.*, 89(3):38009, 2010.
- [185] Owen T Courtney and Ginestra Bianconi. Generalized network structures: The configuration model and the canonical ensemble of simplicial complexes. *Phys. Rev. E*, 93(6):062311, 2016.
- [186] Kiriil Kovalenko, Irene Sendiña-Nadal, Nagi Khalil, Alex Dainiak, Daniil Musatov, Andrei M Raigorodskii, Karin Alfaro-Bittner, Baruch Barzel, and Stefano Boccaletti. Growing scale-free simplices. *Commun. Phys.*, 4(1):1–9, 2021.
- [187] Abigail Hickok, Yacoub Kureh, Heather Z Brooks, Michelle Feng, and Mason A Porter. A bounded-confidence model of opinion dynamics on hypergraphs. *arXiv preprint arXiv:2102.06825*, 2021.

CEU eTD Collection