

UNIVERSITÀ DEGLI STUDI DI TORINO
SCUOLA DI SCIENZE DELLA NATURA
Corso di Laurea Magistrale in Fisica dei sistemi complessi



**UNIVERSITÀ
DI TORINO**

Tesi di Laurea Magistrale

**Understanding cooperative behavior
in online games:
a partial entropy approach**

Relatore:
Prof. Giovanni Petri

Correlatore:
Dr. Michele Starnini (UPC)

Candidato:
Andrea Pio Maria Cota
866282

Controrelatore:
Prof. Piero Fariselli

Anno Accademico 2022/2023

To my family,
they have always known.

Abstract

In social phenomena the basic constituents are humans, i.e., complex individuals who interact and cooperate with a limited number of peers compared to the total number of people in the system, with the aim of achieving a common goal. Understanding and modeling collaborative phenomena is important to explain and predict what really happens in the society we live in, in terms of culture, laws, social contexts and much more.

More recently, we have also begun to interact with others in an online setting, for example we can share our personal updates on social networking sites, engage with and comment on blog entries, but also participate in virtual communities, such as the case of "Reddit Place", a social experiment created by one of the most famous social networks, Reddit, in 2017.

Participants (users) had access to a drawing canvas where they could change the color of one pixel at every fixed time interval. Users were not grouped in teams nor were given any specific goals, yet they organized themselves into a cohesive social structure and collaborated to the creation of a multitude of artworks.

The advantage of online social interactions is that they can be easily recorded in the form of dataset and then studied for scientific purposes benefiting research, in fact, in this work we performed a study on a dataset comprising more than 16M user actions, recorded on the online social experiment, with the goal of discovering how people collaborate and organize on online social games.

Since collaboration patterns are difficult to capture when the relationships between actors are not directly observable, in this thesis we are going to use information theory tools in order to identify specific patterns of collaboration; one of these is Partial Entropy Decomposition (PED) which allows us to distinguish the behavior of users who cooperate from those who do not.

To understand the community dynamics, we, made several definitions of interaction between users, some based only on the placement of their pixels others that also included their color. Moreover, thanks to the timelapse of the experiment, we noticed how different artworks had different behaviours during the whole time: some of these were disappearing, while others remained in the same positions, probably indicating the presence of a conflict, in the former case, and peaceful relations, in the latter.

Finally, we attempted to replicate the dynamics of the canvas on which users interact by generating synthetic timeseries using Gillespie's algorithm and our definitions of interactions, comparing the various results with the entropy measures mentioned earlier showing how PED is an extremely useful tool to identify collaboration patterns also in an online context.

This work shows how complex can be the dialogue between the structure of a system and its emergent cooperative behavior but also offers some interesting insights especially about the use of PED in an online high-order interactive environment, such as r\Place.

Italian abstract

Nei fenomeni sociali i costituenti di base sono gli esseri umani, cioè individui complessi che interagiscono e cooperano con un numero limitato di pari rispetto al numero totale di persone nel sistema, con lo scopo di raggiungere un obiettivo comune. Comprendere e modellare i fenomeni collaborativi è importante per spiegare e prevedere ciò che accade realmente nella società in cui viviamo, in termini di cultura, leggi, contesti sociali e molto altro.

Più recentemente, abbiamo anche iniziato a interagire con gli altri in un ambiente online, ad esempio possiamo condividere i nostri aggiornamenti personali sui siti di social networking, impegnarci e commentare i post dei blog, ma anche partecipare a comunità virtuali, come il caso di "Reddit Place", un esperimento sociale creato da uno dei più famosi social network, Reddit, nel 2017.

I partecipanti (utenti) avevano accesso a una tela da disegno dove potevano cambiare il colore di un pixel a ogni intervallo di tempo prefissato. Gli utenti non erano raggruppati in squadre e non avevano obiettivi specifici, tuttavia si sono organizzati in una struttura sociale coesa e hanno collaborato alla creazione di una moltitudine di artwork.

Il vantaggio delle interazioni sociali online è che possono essere facilmente registrate sotto forma di dataset e poi studiate per scopi scientifici a beneficio della ricerca; infatti, in questo lavoro abbiamo effettuato uno studio su un set di dati comprendente più di 16M di azioni degli utenti, registrate nell'esperimento sociale online, con l'obiettivo di scoprire come le persone collaborano e si organizzano nei giochi sociali online.

Poiché i modelli di collaborazione sono difficili da catturare quando le relazioni tra gli attori non sono direttamente osservabili, in questa tesi utilizzeremo strumenti di teoria dell'informazione per identificare specifici modelli di collaborazione; uno di questi è la Partial Entropy Decomposition (PED) che ci permette di distinguere il comportamento degli utenti che collaborano da quelli che non lo fanno.

Per comprendere le dinamiche della comunità, abbiamo realizzato diverse definizioni di interazione tra gli utenti, alcune basate solo sul posizionamento dei loro pixel altre che includevano anche il loro colore. Inoltre, grazie al timelapse dell'esperimento, abbiamo notato come i diversi artwork avessero comportamenti diversi durante tutto il tempo: alcuni sparivano, mentre altri rimanevano nelle stesse posizioni, indicando probabilmente la presenza di un conflitto, nel primo caso, e di relazioni pacifiche, nel secondo.

Infine, abbiamo cercato di replicare la dinamica della tela su cui gli utenti interagiscono generando serie temporali sintetiche utilizzando l'algoritmo di Gillespie e le nostre definizioni di interazioni, confrontando i vari risultati con le misure di entropia

citate in precedenza e dimostrando come la PED sia uno strumento estremamente utile per identificare i pattern di collaborazione anche in un contesto online.

Questo lavoro mostra quanto possa essere complesso il dialogo tra la struttura di un sistema e il suo comportamento cooperativo emergente, ma offre anche alcuni risultati interessanti, soprattutto per quanto riguarda l'uso della PED in un ambiente online di interazioni high-order, come r\Place.

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 1 |
| 2 | Related works | 4 |
| 2.1 | Social systems | 4 |
| 2.2 | Complex systems cooperation tools | 5 |
| 2.3 | Cooperation | 6 |
| 2.4 | Previous studies on the subreddit | 7 |
| 2.4.1 | r\Place | 8 |
| 3 | The "Place" dataset | 11 |
| 3.1 | Reddit | 11 |
| 3.2 | Place | 12 |
| 3.3 | Dataset's structure | 14 |
| 4 | Methods | 16 |
| 4.1 | How to identify and select artworks | 16 |
| 4.2 | Defining "interactions" | 18 |
| 4.2.1 | First definition of interaction | 19 |
| 4.2.2 | Null model | 19 |
| 4.2.3 | Different timedeltas | 20 |
| 4.3 | Understanding the dynamics | 20 |
| 4.3.1 | GilleSpie Algorithm | 20 |
| 4.3.2 | Generating synthetic timeseries | 22 |
| 4.4 | Partial Information Decomposition | 23 |
| 4.4.1 | PID intuition | 24 |
| 4.4.2 | Partial Entropy Decomposition | 26 |
| 4.4.3 | Our PED configurations | 29 |
| 5 | Results | 31 |
| 5.1 | Exploratory data analysis | 31 |
| 5.1.1 | Activity trend analysis | 31 |
| 5.1.2 | Inter-event times analysis | 33 |
| 5.1.3 | Activity through the time | 34 |
| 5.1.4 | Color feature analysis | 37 |
| 5.2 | Identifying collaborations and conflicts | 38 |
| 5.3 | Interactions analysis | 43 |

CONTENTS

| | |
|---|-----------|
| 5.4 PED analysis | 46 |
| 5.4.1 Inside & outside the artworks | 48 |
| 5.4.2 Toy models | 55 |
| 6 Discussions | 62 |
| 6.1 Future works | 64 |
| Bibliography | 68 |
| Ringraziamenti | 71 |

Chapter 1

Introduction

Social research helps to acquire knowledge of human society. The dynamic nature of society makes it challenging to understand to move towards progress and welfare. Social research enables us to analyze social behavior, understand the causes, and accelerate its evolution.

Exploring social phenomena is extremely important due to its profound impact on our society and individual lives. It provides us useful information about the intricate world of human interactions and behaviors. By delving into the dynamics of how people engage in various social settings, we gain a deeper understanding of their behavioral norms and customs.

When we examine social phenomena, we can identify trends, anticipate potential issues, and formulate informed strategies to tackle them. Information gathered from the research can help us predict the behaviour of certain individuals or groups. When we have a good understanding of a social phenomenon, we may have a better idea of how to govern or guide it.

The data collected in these contexts can be used to bring planned structural changes to social life. It can help us understand the current needs people have and develop an action plan to meet those needs.

With the raise of social media, we have at our disposal large amounts of data and digital traces of social activities that could help to improve the comprehension of large-scale social phenomena through a quantitative background ranging from mobility patterns [1], to information spreading [2].

Consequently, for the purposes mentioned above, the presence of this data can be of great help; however, there are still a number of limitations: social media data is often unstructured, noisy, incomplete, and biased. For example, some users may post fake or misleading information, some may delete or edit their posts, and some may use different accounts or platforms. Moreover, social media data is constantly changing and evolving, which makes it difficult to track and compare over time. Also, social media data involves personal and sensitive information of users, such as their opinions, preferences, behaviors, and identities. Finally, social media data is often complex, multidimensional, and contextual. It requires a combination of quantitative and qualitative skills to analyze and understand the data. For these reasons, analysts

CHAPTER 1. INTRODUCTION

need to use appropriate techniques and tools to visualize, summarize, and present the data accurately.

Once one is aware of what it means to use social media data, it is also appropriate to talk about emergence, which is one of the most important features of complex systems. Emergent phenomena result from the interactions of individual entities. By definition, they cannot be reduced to the system's parts: the whole is more than the sum of its parts because of the interactions between the parts. An emergent phenomenon can have properties that are decoupled from the properties of the part. For example, a traffic jam, which results from the behavior of interactions between individual vehicle drivers, may be moving in the direction opposite that of the cars that cause it. This characteristic of emergent phenomena makes them difficult to understand and predict: emergent phenomena can be counterintuitive.

Agent Based Models (ABM) [3] is, by its very nature, the canonical approach to modeling emergent phenomena: in ABM, one models and simulates the behavior of the system's constituent units (the agents) and their interactions, capturing emergence from the bottom up when the simulation is run. It is an approach to simulate the behavior of a complex system in which agents interact with each other and with their environment using simple local rules.

Besides experimental approaches, one can represent emergent phenomena using theoretical approaches, like games theory, i.e., the study of mathematical models of strategic interactions among rational agents. It has applications in all fields of social science, as well as in logic, systems science and computer science. Among others, games theory is a tool that first defined the concept of "cooperative game" and non-cooperative game, using mathematical definitions such as the payoff matrix. Games theory is useful for all those situations in which users interact in a way that is not well defined, or simply not known in advance. For example, one can have the situations where people cooperate, or compete, or even in which the competitors try to both positively influence their own returns, while negatively affecting those of their competitors [4].

That said, in this work the case-study system is the canvas of Reddit Place 2017 edition: the event organizers had a surprising plan for the Reddit community on April Fool's Day. They kept the proposed activity a secret, known only to the developers until the event started. The participants were introduced to an empty canvas with a brief message: "You can place a tile on it, but you must wait for your next turn. Alone, you can create something, but together, you can create something greater."

There were no additional instructions except for a link to access the graphical interface, which provided a large 1000x1000-pixel blank canvas. Each participant had the freedom to select a tile and color it from a 16-color palette. However, they had to wait at least five minutes before making another interaction. This time constraint discouraged individual actions, promoting collaborative efforts.

The experiment lasted for three consecutive days and attracted over one million users who made more than 16 million interactions. Although Reddit did not specify any purpose or task, users quickly organized themselves into communities that collaborated

CHAPTER 1. INTRODUCTION

to create and maintain images representing their identities, such as national flags, logos, video game characters, manga, sport team crests, famous paintings, and other iconic images from online culture. As the days progressed, the canvas transformed from chaos to a harmonious mosaic of pixel art, with images bordering and intertwining with each other, forming a seamless and visually stunning artwork.

Anyway, in organized complex systems like Place, composed of large numbers of parts with non-trivial interactions, the only way forward is to tackle the complexity head on, using computer simulations, structure learning, and large models. This way forward is "information theory," which is one of the most used tool in this field of research [5].

Information theory is a branch of applied mathematics and computer science that deals with the quantification, storage, transmission, and processing of information. At its core, information theory aims to measure and understand the amount of information contained in a message, signal, or data set. It introduces the concept of "entropy" as a measure of uncertainty or randomness within a given set of data. In this work we are going to use "Partial Entropy Decomposition" (PED) [6]: it provides a framework with which we can extract all of the meaningful structure in a system of interacting random variables.

This thesis work will show how users organized their actions in limited spatio-temporal areas and how PED is an incredibly useful tool for detecting patterns of collaboration between users even in online game environments by identifying the different type of information contained in collaborative and non-collaborative relationships.

Chapter 2

Related works

In this section, we will delve into the existing body of literature pertinent to our case study. We'll begin by exploring broader concepts such as social systems, and cooperation. Subsequently, we will progress towards research that has delved into the Place experiment, gradually narrowing our focus to the specifics of our study.

2.1 Social systems

A social system essentially consists of two or more individuals who interact directly or indirectly under limited circumstances. Physical or territorial boundaries may exist, but the basic sociological reference point is that individuals can collaborate, i.e., align toward a common focus. Therefore, it is appropriate to consider different kinds of relationships, such as small groups, political parties, and society as a whole, as social systems. Also, social systems are open systems that exchange information with and are often linked to other systems.

Examples of social systems are abundant in various fields: nations (societies that hold territories with a formal government), economic systems (they produce and distribute value by putting capital to work), cities (places where people live in organized with a local government), media (the internet has allowed anyone to publish information to the world using platforms such as blog or social media) and many others.

Collaboration within social environments is a crucial aspect also of complex systems. In such settings, individuals interact, share information, and collectively contribute to the development of shared goals or outcomes. Cooperation can give rise to emergent behaviors that shape the overall dynamics of a social system. For instance, in a community, individuals might collaborate to address common challenges, leading to the emergence of collective solutions that benefit the entire group. However, the balance between individual interests and collective goals can influence the success and sustainability of collaboration within complex social systems.

Understanding and managing social systems, particularly in the context of collaboration, is essential for addressing various real-world challenges. For these reasons, and also as this thesis work also wants to contribute, researchers in the fields of science, like sociology, studied different aspects of social systems, varying from communication

between groups of animals [7], to information spreading on online platforms [2], like Reddit.

2.2 Complex systems cooperation tools

When talking about complex systems, it is necessary to understand that they can encompass several scientific fields and different domains, including the social field as is the case of this thesis.

Consequently, the tools needed to analyze them can also be of various kinds, and thus arise from different contexts, such as computer science, physics, economics and more.

For example, concepts from Games theory, defined as "the study of mathematical models of strategic interactions among rational agents"¹, are often used to reproduce a specific phenomenon through the creation of rational agents, i.e., those who carry out certain actions with the sole aim of maximizing their gain, and interact with each other and with a given environment. One of the most famous games and one that has inspired others is that of the "prisoner's dilemma", who involves just 2 agents that can cooperate or not.

Moreover, agent-based models (ABM) are computational models very used in the field of complex systems to simulate the behavior and interactions of autonomous agents (individuals and collective entities such as organizations or groups). We have several examples of their applications in different fields like biology [8], or epidemiology [9], and many more.

Moving on, social networks were designed following network theory fundamentals; it is the theory that allows networks to be defined as graphs, that is, structures consisting of vertices and connected through links. As an example, social communities are networks in which the nodes are people (or organizations of people) between whom there are many different types of possible relationships. More linked to cooperation, Cohen, Havlin, and benAvraham [10] propose an intriguing public health strategy that employs "network thinking" for vaccination strategies in the face of limited resources: Instead of vaccinating people randomly, ask a large number of people to each name a friend, and vaccinate that friend. The idea is that in a social network with power-law degree distribution and high clustering coefficient, collecting names of friends of individuals will be more likely to produce hubs of the network (people with many friends) than choosing random individuals to vaccinate. And vaccinating a hub is more likely to slow the spread of a disease than is vaccinating a randomly chosen person.

Finally, as in this work, information theory is one of the most used frameworks in this field of research. It is the mathematical treatment of the concepts, parameters, and rules that govern the transmission of messages through communication systems². In [6] Varley et al. find that synergistic information (i.e., cannot be reduced to pairs of

¹from R. B. Myerson, Game theory: analysis of conflict. Harvard university press, 1991

²from: ScienceDirect

nodes), is invisible to standard network-based approaches. Their results provide strong evidence that there exists a large space of unexplored structures in human brain data that have been largely missed by a focus on bivariate network connectivity models. This, and other related works, provide a very general approach for understanding higher-order structures in a variety of complex systems.

2.3 Cooperation

The emergence of cooperation is a fundamental feature of complex systems. This phenomenon resonates in diverse domains, ranging from the microcosm of microbiology [11] to various possible social systems. This natural inclination to cooperation takes different forms, requiring exploration through as many different approaches.

One of the most important paradigms for delving into cooperative dynamics is game theory. The emblematic example is the Prisoner’s Dilemma, a theoretic concept in which two rational agents must choose between cooperation and conflict. The resulting payoff, i.e., their reward, depends on the synchronized choices of these agents. This paradox carries a fundamental challenge: the act of cooperation often carries costs, while exaggerated self-interest can tip the scales in favor of conflict. The fundamental question remains: why does cooperation materialize?

A broader perspective reveals that such concretization of cooperation lies in its potential to generate new (in some cases more convenient) payoff levels that would be unattainable by individual efforts alone. This fundamental principle is well established, but the mechanisms and circumstances that generate it are still the subject of extensive investigation.

In 1981, Axelrod revolutionized this concept by presenting a seminal study that sought to explain the riddle of the Repeated Prisoner’s Dilemma [12]. Among the strategies, the "TIT FOR TAT" emerged as one of the most convincing strategies. Initially, the player, or agent, cooperates and later repeats the opponent’s previous move. There are many parameters that influence agent cooperation, for example, the initial prevalence of defectors, or even stochastic errors or the influence of time horizon on decision dynamics.

Recently there has been a migration of evolutionary games into the domains of graphs and hypergraphs [13]. This transition reflects the evolving attempt to examine the interaction of cooperation within network structures. Although the precise interrelationship between network topology and the establishment of cooperation in agent-based models remains an enigma, certain attributes of real-world networks seem to promote inclinations toward cooperation [14]. In particular, the grouping of entities into cohesive communities creates fertile ground for cooperative efforts, while central hubs serve as keystones of cooperation, connecting a myriad of nodes [15].

Moreover, a paradigm shift occurs when dynamic networks are contemplated, in which agents can change their memberships. In this context, cooperation proliferates under specific temporal conditions. Specifically, when interaction frequencies align

with regular temporal patterns, cooperation finds fertile ground. In contrast, the irregular bursts that characterize real-world temporal networks can discourage cooperative outcomes [16].

In general, collaborative networks realize different scenarios. An important example involves the network of scientific collaborations, where connections indicate authorship of a collaborative article [17]. This context exhibits properties such as the "small-world" arrangement where randomly chosen scientists are separated by mere knowing people and a high coefficient of clustering, indicative of entangled collaborations. Remarkably, these networks can be elegantly transformed into hypergraphs by encapsulating interactions within structures that unite co-authors of shared publications.

Finally, to get to more recent times, the advent of online platforms has given rise to new possibilities for collective interaction and collaboration. On all, certainly social networks have emphasized this trend. Beyond the better-known Facebook or Instagram, platforms such as Twitch, which allows unlimited live streaming guaranteeing endless entertainment, and Reddit, which organized the world of forums and comments by establishing specific "threads" where people discuss topics of all kinds, have almost eliminated the barrier of a screen separating two people, strengthening interactions between individuals. Notable experiments, such as Twitch Plays Pokémon³ and "Reddit r/place", further emphasize the presence of these cooperative mechanisms.

In conclusion, we have illustrated many times that the world of complex systems has applications in extremely varied contexts. Insights from game theory, evolutionary game frameworks and network science highlight the mechanisms underlying cooperative behavior. Whether manifested through any of these concepts, cooperation exists, and it is at the basis of emerging models describing networks of interactions.

2.4 Previous studies on the subreddit

As a social network, Reddit has been one of the most studied by researchers especially because the datasets generated by the activity of the users on this platform are often openly available for researchers, in contrast to other sources of data in computational social science.

Despite its recognized quality, one should be careful while using the dataset for research purposes. The risks of missampled data obviously cause commenting or posting rate distortions, missing information in user time series and possible inability to reconstruct certain discussion trees.

That said, we are going to illustrate some of the multitude of works that involved Reddit datasets; the majority of them is related to the fields of computational social science and machine learning.

³is a social experiment and channel on the video game live streaming website Twitch, consisting of a crowdsourced attempt to play Nintendo's Pokémon video games by parsing commands sent by users through the channel's chat room.

A very common phenomenon which frequently happens in online discussions is *trolling*, offensive or menacing messaging. Mojica [18] collected and studied an annotated dataset of trolling comments in discussions on Reddit using a variety of language features.

Kumar et al.[19] considered interactions between communities in the form of mobilization by users of a community (the source of the ‘attack’) for hateful comments on posts from another community (the target of the attack). The authors propose an LSTM (Long short-term memory) neural network model that uses textual and social features in order to identify whether a given post will produce a mobilization.

Recurrent neural networks (RNN) have also been employed to measure community endorsement. Fang et al. [20] constructed an RNN trained to predict comment scores.

Horne and Adali [21] studied how posting news articles on subreddit */r/world-news* influences their popularity and conclude that changing the article titles results in greater popularity comparing to leaving the original one.

Tan and Lee [22] studied the posts of users across the subreddits and found that users on average tend to explore and continuously post in new communities, moreover they tend over time to share their activity evenly between a small number of communities with diverse interests. Differences in posting patterns of users may be used for prediction of the user’s future settlement status in a community.

Finally, Medvedev et al. [23] created a survey that lists a series of Reddit-related works, who also illustrates the richness of datasets that comes from to online discussion platforms, with Reddit as a dominant example.

2.4.1 r\Place

Only a limited number of research papers have delved into this particular experiment. What’s intriguing is that each of these works offers its own unique perspective on the event, revealing distinct aspects of it.

In 2018, scholars Thomas F. Muller and James Winters [24] put forth an analysis using the lens of cultural evolutionary theory. In their approach, they considered artworks as forming lineages of descent, where colored pixels represent cultural attributes. They then use these observations to update their own beliefs and make informed choices about which colored pixel to select and where to position it on the canvas.

This information exchange eventually guides users to adopt certain patterns in their placements, consequently affecting the resulting images. These patterns can be measured using the principles of information theory, employing concepts like Shannon entropy [25] (both globally and locally), as well as image compression. Moreover, specific regions on the canvas (referred to as artworks) maintained low entropy, implying order, while overall entropy remained high. This suggests that diversity was confirmed at a global level, yet organized into predictable distributions at the local level.

During that same year, Ben Armstrong [26] undertook an exploration of this data in his master’s thesis, employing a peer production perspective. He compared the event

CHAPTER 2. RELATED WORKS

to other platforms characterized by peer production, such as Wikipedia. By "peer production", he alluded to the online process through which numerous individuals collaborate to generate information or culture.

Armstrong constructed a model of the experiment utilizing an agent-based approach, investigating two forms of agent coordination. The first is spontaneous coordination, which denotes a user's capacity to engage in collective behavior without centralized planning. The second is external coordination, which entails designing specific regions on the canvas outside the experiment. This design emulates scenarios like communication via subreddits during the event.

Once again, in 2018, Rappaz et al.[27] introduced an approach to distinguish collaboration patterns utilizing a predictive model founded on an embedding technique. In this method, each user is linked to a vector within an embedding space. The coordinates of this vector are continually adjusted based on the user's actions on the canvas. This adjustment follows a machine learning training algorithm that employs proximity in the latent space to approximate social similarities. Essentially, users are more inclined to engage in activities near other participants who are proximate in the latent space, as opposed to those who are distant.

In 2020, Elías Gabriel Gil with his master thesis [28] tried to characterize the online context of participation shared by the users, and to identify the forms of organization established among them for their actions in Place. This work shows up how the users organize their actions through the subreddits, and also reveals some relationships of collaboration and conflicts between communities belonging to different artworks.

Prateek Vachher and colleagues [29] took on the task in 2020 of examining this dataset from the perspective of conflicts over space among communities. They established specific criteria to pinpoint areas on the canvas that experienced intense contention among two or more groups of users. This allowed them to associate communities with these conflict zones and determine winners and losers in such scenarios.

Firstly, they concluded that, in the context of winning conflicts, the number of active users carries more significance than the level of individual user activity. Secondly, their analysis indicated that conflicts often involve multiple communities.

In 2021, Kristina T. Litherland and Anders I. Mørch [30] undertook an examination of the Place event by studying the progression of two distinct types of objects: visual artifacts and social artifacts. In this context, "visual artifacts" encompassed the actions and contributions made on the canvas, while "social artifacts" referred to the interactions taking place through subreddits.

Both types of artifacts were further categorized into two key aspects. Firstly, there was the concept of "bottom-up emergence," which emerged from the collaborative efforts and interactions among users. Secondly, there were "top-down instructions" or rules that were imposed by developers. These rules served as regulatory and control mechanisms for the experiment.

Their findings revealed intricate and compelling interactions between these two types of artifacts. These interactions played a crucial role in facilitating the creation,

CHAPTER 2. RELATED WORKS

stabilization, and preservation of a specific artwork (the Mona Lisa image) within the Place event.

Also, in 2022, with his master thesis [31] Liber Dorizzi studied the relationship between the structure of a real social system (people discussing on subreddits) and its emergent behavior (same people who generated the final canvas). To study the activity of the users he used an embedding space in order to build a predictive model on individual user actions. The embedding proved to be a very appropriate tool for modeling the latent structure of users' decision-making process.

As a result he found some interesting relationships between structure and activity, for example, generally groups that created artworks had a more cohesive social structure than uniform samples and also slightly greater than users within subreddits that did not generate artworks.

In 2022, Alyssa M. Adams and colleagues [32] conducted research aimed at addressing a specific question: what guiding principles do users adhere to when placing a new pixel on the canvas, considering the current state of the canvas itself? This study deliberately excluded social communication through subreddits and focused exclusively on user engagement with the canvas. Firstly, they utilized a statistical rule-based model, which represented actions on the canvas akin to a 2D cellular automaton. Secondly, they employed a CNN-based (convolutional neural network) neural network approach. The researchers identified that the rules dictating the emergence of objects on the canvas predominantly stem from external social collaborations that occur within subreddits. This insight emphasizes the intricate interplay between external social dynamics and inherent canvas interactions in shaping the evolution of the Place event.

Chapter 3

The "Place" dataset

In the present era, social media platforms present a terrific opportunity for studying extensive datasets concerning genuine human interactions and social systems. These platforms are extensively utilized to investigate various aspects and phenomena of human society, encompassing sentiment analysis, opinion dynamics, the spread of rumors, the echo chamber effect, detection of fake news, peer production, online gaming cooperation, and more.

Researchers make an assumption when utilizing datasets derived from online interactions, considering that virtual behavior can serve as a reliable proxy for studying real-life behaviors. In this study, we focus on analyzing a specific dataset sourced from the social network Reddit, which ranks among the largest online social platforms available today.

3.1 Reddit

Reddit is a popular social media platform and online community that allows users to share and discuss content on various topics. It operates as a network of user-generated communities, each centered around specific interests, known as "subreddits." These subreddits cover a wide range of subjects, including technology, science, gaming, news, art, sports, and much more.

The dynamic by which users act on the social network is very simple: users can subscribe to subreddits based on their interests. Each subreddit has its own theme and rules, and members can submit posts and engage in discussions related to the topic.

In this context a "voting system" is involved to determine the visibility of content. Posts can be upvoted if they are considered valuable or interesting, and downvoted if they are deemed irrelevant or inappropriate. The number of upvotes and downvotes determines a post's position on the subreddit's page.

Obviously, we have also the presence of comments and discussions in which users are engaged with other community members. The most upvoted comments usually appear at the top of the comment section and each subreddit has its own set of moderators who enforce the rules and guidelines of the community. They have the authority to remove posts or ban users who violate the rules.

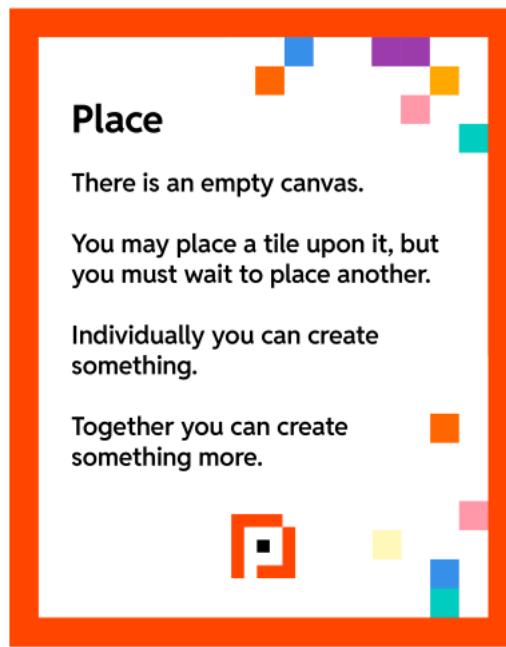


Figure 3.1: Reddit place description

Some of the latest features added are the Reddit "Gold" version and the Awards: users can purchase "Reddit Gold" (now called "Reddit Premium") to access additional features and to support the platform. Additionally, users can give awards to posts and comments as a form of recognition and appreciation.

Reddit has grown into one of the largest and most influential social media platforms, allowing users to discover and engage with content, share their opinions, and participate in diverse communities centered around their interests. The platform's dynamic nature and community-driven content curation have contributed to its widespread popularity and ongoing success.

3.2 Place

On April Fools' day in 2017, the social media platform Reddit introduced a fascinating social experiment called "Place." The experiment revolved around an initially empty canvas consisting of a 1000x1000 grid of tiles (or pixels). All Reddit users who had registered their accounts before the event were welcome to actively participate. They could access the interface, pick a tile, and modify its color using a 16-color palette. Users had the freedom to either color an empty tile or overwrite an already colored pixel. We'll refer to the latter action as "placement."

To ensure a collaborative atmosphere and prevent individual dominance, users faced a limitation on their action rate. After making a placement, they had to wait for a fixed duration before their next action on the canvas, typically 5 minutes (20 minutes for unverified accounts). During this waiting period, any other regular Reddit user

CHAPTER 3. THE "PLACE" DATASET

could recolor their tile.

The main objective behind this rule was to discourage isolated actions and encourage users to coordinate with others, aiming to create meaningful patterns and collective artworks on the canvas.

The experiment lasted for a continuous period of 72 hours and garnered significant engagement and participation from the Reddit community. Over one million unique users actively took part, resulting in a remarkable total of more than 16 million interactions on the canvas.

The experiment's objective was open-ended, with no specific goal or preliminary task provided, except for the concise instructions given in the developers' opening post:

"There is an empty canvas. You may place a tile upon it, but you must wait to place another. Individually you can create something. Together you can create something more."

In addition to the developers' opening post, there were only four general rules for the "r/place" experiment: be creative, maintain civility, follow site-wide rules, and refrain from sharing personal information. Beyond these rules, no further instructions were given.

Hence, the final outcome of the experiment was unpredictable, but it quickly transitioned from a chaotic state to a state of self-organization within the first day. Users naturally formed groups and collaboratively created meaningful figures on the canvas. As the second day began, the blank tiles on the canvas were all taken, leading to interactions and negotiations between the groups for available space. This interaction sometimes resulted in conflicts as different groups competed for control.

By the end of the third day, the canvas had transformed into an iconic mosaic of art, where countless works have been entangled together, reflecting the collaborative efforts and conflicts of hundreds of thousands of users. The final canvas showcased recognizable themes such as national flags, video game logos, comics, sports teams, and various other subjects from online culture. This observation led to the hypothesis that these themes might be linked to specific subreddits, indicating a potential connection between subreddit structures and the activities seen on the canvas.

Further evidence of the strong connection between subreddits and artworks emerged when a collaborative effort produced an atlas a few weeks after the experiment's conclusion. This atlas documented all the artworks featured in the final canvas, providing titles, brief descriptions, and, where possible, the main subreddits that contributed to their creation. This valuable resource, containing around 1500 artworks, serves as a ground truth for part of our future analyses.

Technically, the dataset from the experiment is readily accessible online and organized in a dataframe format. Each row in the dataframe represents one interaction (placement) on the canvas, resulting in an extensive collection of over 16.5 million rows. By defining specific spatiotemporal boundaries on the dataframe, we can observe and analyze the flow of interactions in selected regions of the canvas and within particular time intervals.

Also, a time-lapse video of the entire experiment is available on YouTube at the fol-

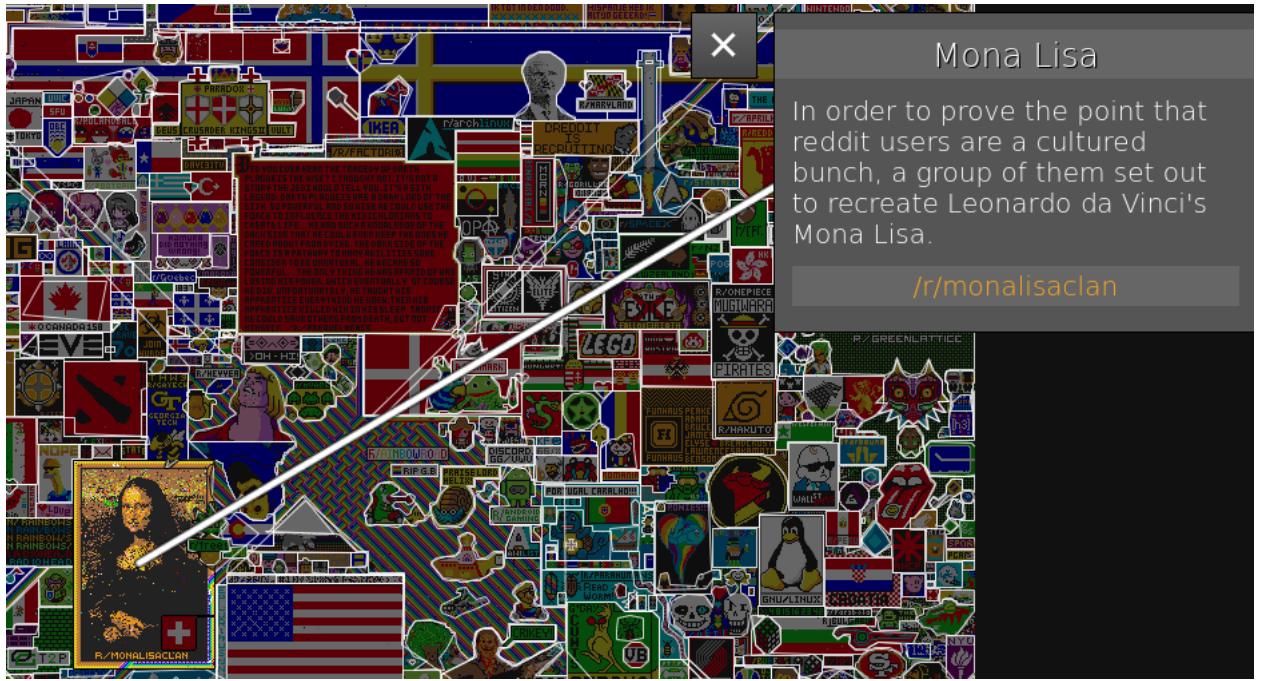


Figure 3.2: Reddit place artworks atlas

lowing link: <https://www.youtube.com/watch?v=XnRCZK3KjUY>. This video provides a visual representation of how the canvas evolved over time, capturing the collaborative and dynamic nature of the project.

3.3 Dataset's structure

The dataset of Reddit Place 2017 is a collection of pixel data, user interactions, and metadata that offers a comprehensive view of how a massive collaborative artwork emerged over time through the contributions of Reddit users. It provides a unique view into online collaboration, creativity, and the dynamics of a virtual community.

In detail we have a data frame of 16 millions rows, each one representing a placement, and 5 columns, that are:

- "ts" (int type): it's the timestamp of the placement;
- "user hash" (str type): the hashed username of the user that placed that pixel;
- "x coordinate" (int type): x coordinate on the canvas of that pixel;
- "y coordinate" (int type): y coordinate on the canvas of that pixel;
- "color" (int type): the specific color of the placed pixel on a scale of 16 different colors.

In this way the colored canvas has been created (3.3).

CHAPTER 3. THE "PLACE" DATASET



Figure 3.3: Some of the artworks of the final canvas

Chapter 4

Methods

In this chapter we will review the methods used to analyze the dataset. We divided the various methods according to the following progression:

- First, we explain how to identify and study an artwork;
- Then we defined our concept of interaction between users based on space and time, to identify some relationships between them;
- After that, with the aim to understand how users behaved in the canvas, we generated some synthetic timeseries using GilleSpie algorithm;
- Finally, we performed Partial Entropy Decomposition on specific artworks to show the presence of collaboration patterns.

4.1 How to identify and select artworks

As previously stated, artworks naturally appear to be linked with specific themes, and as a result, they often have connections to particular subreddits.

We view the dataset as an extensive space-time container, where each artwork placement occupies a distinct position within it.

Consequently, our objective is to pinpoint the origin of an artwork within both time and space dimensions. Additionally, we aim to identify the cohort of users responsible for its creation.

To accomplish this, we adhere to the following procedure:

- Utilizing either timelapse recordings or screenshots of the event, we pinpoint an artwork for analysis, along with the specific pixel coordinates defining its location.
- We make a rough determination of two critical moments: the point when the artwork is completed and, in reverse order, when it starts to emerge on the canvas (fig. 4.2).

Placements = 10438
Users involved = 3720

Detail - $x \in (31, 85)$ $y \in (713, 752)$ pixels

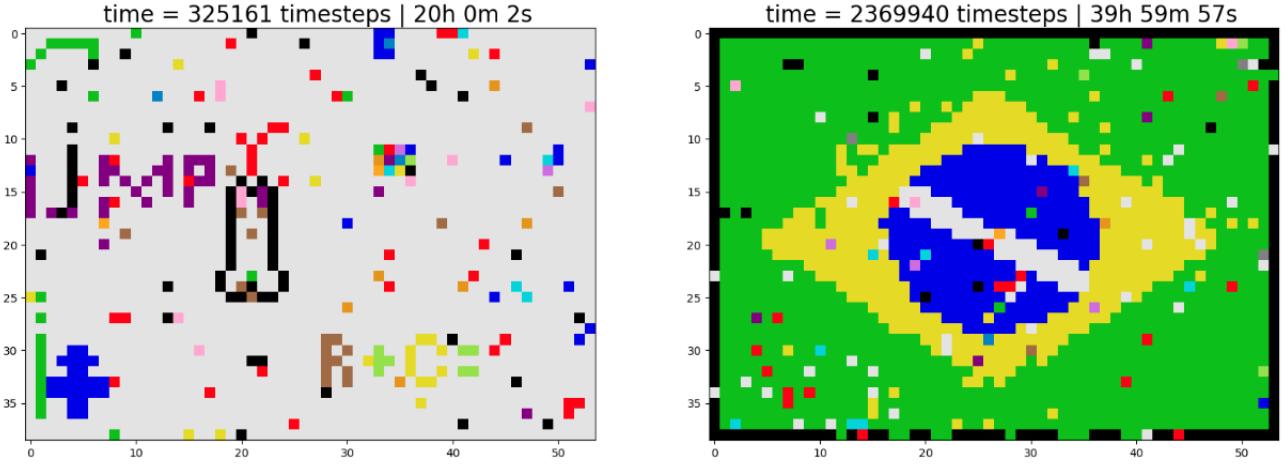


Figure 4.1: An example: brazil artwork at the beginning and at the end.

- Once we've established this spatiotemporal framework, we proceed to select the users who engaged within this timeframe.

The information we gain by doing this are visible in (fig. 4.1): number of placements and number of users involved in the space-time area, x and y coordinates and also the start and the end timestamps.

We have to point out a few things: of course our conception of beginning and end of an artwork is biased, because nobody can be sure about the specific moment of birth of the artwork, but approximately we can try to identify it. This consideration doesn't change anything in our analysis, we are not going to study the behaviour of every user in every artwork, but we just want to have a panoramic view of what it's happening in the area we choose to study.

Artworks can assume various forms, ranging from highly regular to more irregular shapes. They may be localized or extend across a significant spatial area. For instance, consider the "rainbowroad," a rainbow-colored stripe that nearly traverses the entire canvas. Conversely, some artworks lack any defined shape, such as the "amorphous void" or the "blue corner." Irregularly shaped artworks pose challenges in terms of framing.

Consequently, our attention is primarily directed toward more regularly shaped artworks, such as rectangular designs for flags and logos, which are easier to confine within a specific timeframe and spatial context.

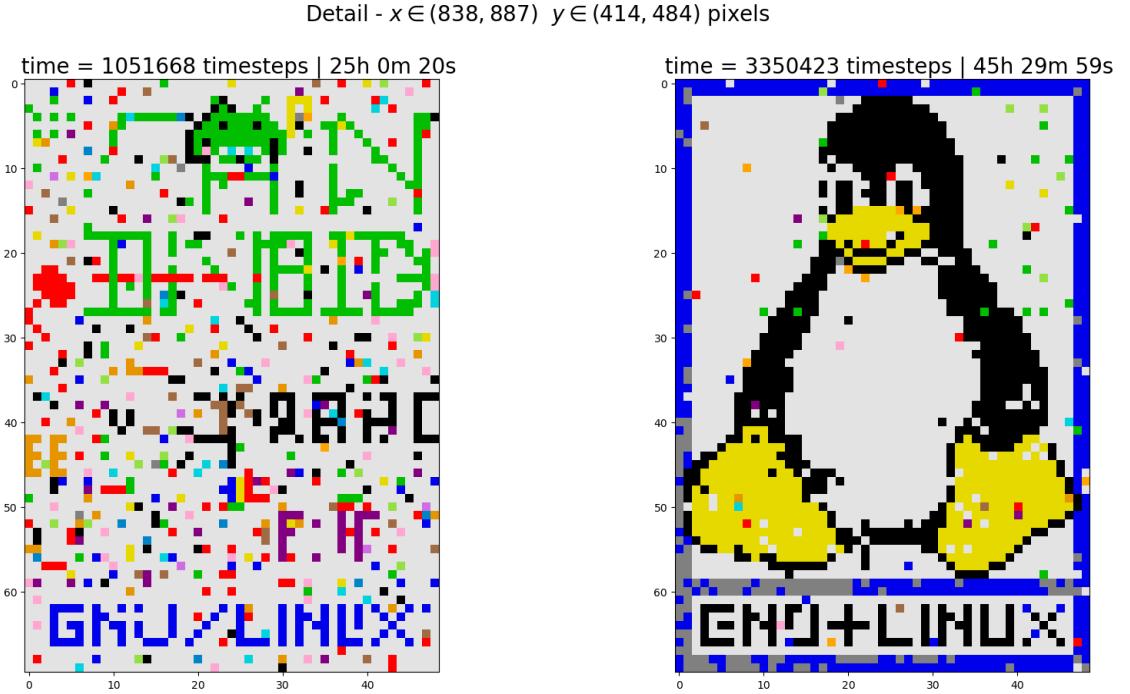


Figure 4.2: The moment when an artwork begins and ends, in this case Linux artwork

We focused our analysis on 16 different artworks with regular shape: *argentina*, *suomi*, *brazil*, *germany*, *ue*, *linux*, *onepiece*, *vangogh*, *starwars*, *suomi2*, *mona*, *blackvoid1*, *blackvoid2*, *naruto*, *dota2*, *rick*.

4.2 Defining "interactions"

We know for sure that people on our environment has *interacted*. The main question is: how they did that?

How described in [28], people in Reddit is organized into communities, that, in the context of Reddit Place, simply are *subreddits*. The cited article describes what happened only in the case of the "argentina" subreddit, but we can assume that a similar situation has existed in the other subreddits. Basically through discussions on the specific subreddit, users organized their placements very precisely in terms of positions and colors. Also this threads were useful to decide where to collaborate or not with neighbours artworks.

So this what happened outside the canvas, we need now to identify this type of organization in the canvas, clearly this is the hard part of the work.

We start by defining our concept of *interaction* in the canvas.

However, over time, its definition has been modified based on the results.

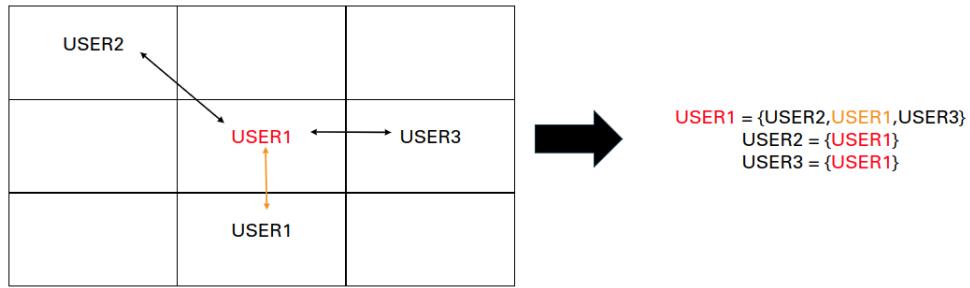


Figure 4.3: Visual representation of our definition of interaction.

4.2.1 First definition of interaction

The main features that could help us into the task of the definition of an interaction on the canvas are: position (in terms of x and y coordinates), timestamp and color.

So we can say that: *an interaction will occur if in the last 10 minutes since the timestamp of a placed pixel by a user, other users or the same user are/is in one of the 8 surrounding positions of the central one.* At the moment we don't include "color" feature in this definition.

In this scenario we could have two different types of interactions: a "self" interaction, when is the same user the one who placed into the central position, or an "other" interactions, that means by two different users. Clearly we could have maximum 1 self interaction in a timedelta of 10 minutes because every user has to wait 5 minute to place the next pixel.

In (4.3) we offer a visual representation of this definition of interaction.

Thanks to this, we can obtain a sparse matrix of the interactions, in which each row represent a single user on the canvas with his number of interactions with others, or with himself on the diagonal (self interactions).

4.2.2 Null model

Our definition of interaction is just a starting point, is obviously not the most accurate measure of users encounters.

The space in which our users live in is limited, it consists of 1 million unique different positions, but also we have more than 1 million unique users, so there is a non null probability that users interacted without intending it. This is what we call *random effect*.

For this reason we need to know how much *random effect* is present in our matrix. To do this we can use a *null model*: it consists in doing a reshuffle of the positions (x and y coordinates) of the original dataset and then calculate the number of interactions present in this new dataframe, that will be stored into a null matrix.

Hopefully the number of interactions present in the null matrix will be very lower than the real matrix, quantifying the random effect.

In this way we just need to subtract the random interactions to the original one to get the "real" number of them.

4.2.3 Different timedeltas

In fig.(4.4) we can see the probability distribution function (PDF) of the inter-event times of every user in the experiment.

What we see is a function that has a peak near 300 seconds. The reason of this is obviously related to the 5 minute bond that every user has to respect to place another pixel.

But why we have non null values before the 300s? Well, the answer is related to the fact that the timedelta of this bond did not remain the same during the whole experiment, instead it's changed in 1 minute, 10 minutes, 20 minutes and so on. This explains this trend.

Now, in the context of our definition of interaction, if we think about a "self" interaction (between the same user) we know for sure that in a timedelta of 10 minutes, like we have defined before, we can have maximum 1 occurrence of this interaction. The important thing to understand is when the majority of them has occurred, if it's just after the 5 minutes, or exactly before the 10 minutes or any other chances. The same goes for the other type of interaction (between two different users), that is, we need to understand when users interact the most in a well defined timedelta.

For these reasons we modified our first definition of interaction by selecting different timedeltas (10 seconds, 30 seconds, 1 minute) but maintaining the rest of the definition as it was.

The findings of this analysis are stored in table (5.2).

4.3 Understanding the dynamics

One of the aims of this work is to understand which reasons motivate a user to put a pixel in a specific place, of that specific color, on that precise time interval. Consequently, model the behaviour of these people could be helpful to discover which are the dominant degrees of freedom of the system under investigation.

In this section, we explain how we tried to replicate the canvas dynamics and which tools helped us in this objective.

4.3.1 GilleSpie Algorithm

As described in [33], several complex systems are guided by interactions among their components through timestamped discrete occurrences. As an illustration, in chemical systems, a chemical reaction event alters the quantity of reagents of specific kinds in a distinct fashion.

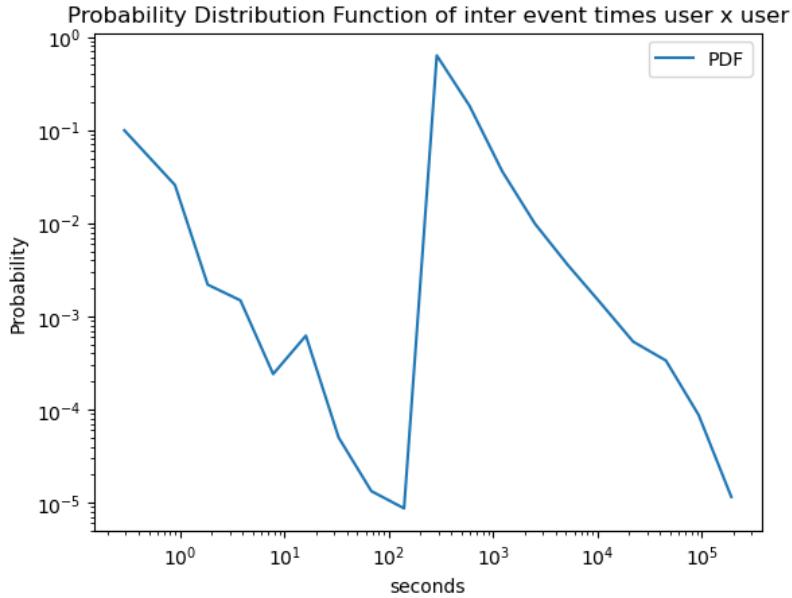


Figure 4.4: PDF of inter event times per user

Stochastic point processes, notably Poisson processes assuming that events unfold independently and at a steady pace over time, serve as a fundamental tool for describing the dynamics of chemical systems. They also prove valuable in the emulation of epidemic phenomena within a population and numerous other systems.

Consider an occurrence-based system in which events arise from Poisson processes running concurrently. In chemical reaction systems, each Poisson process (potentially with distinct rates) is linked to a particular reaction. In the context of epidemic occurrences within human or animal contact networks (that is, graphs), each Poisson process is assigned to an individual or a connection, which could potentially transmit the contagion. The event frequency of certain Poisson processes may shift in response to a reaction or infection occurring within the overall system. The most straightforward simulation approach involves time discretization and the assessment of whether an event takes place within each time interval for individual processes. This extensively employed technique is less than optimal as it demands a small enough time interval to reach high precision, incurring computational costs.

The Gillespie algorithm presents an efficient and statistically precise approach for multivariate Poisson processes. Specifically, the direct method of Gillespie exploits the notion that the summation of independent Poisson processes results in a singular Poisson process, with its event rate being the sum of the constituent Poisson processes. Capitalizing on this mathematical property, the Gillespie algorithm only necessitates the emulation of a single Poisson process.

Here is how does it works: the original algorithm assumes that there are N independent Poisson processes running in parallel, each with a rate of λ_i (where $1 \leq i \leq N$). Because these Poisson processes are independent of each other, their combination results in a new Poisson process with a rate equal to the sum of the individual rates,

which is $\sum_{i=1}^N \lambda_i$. Consequently, the algorithm begins by drawing a time increment, Δt , to the next event of this combined Poisson process from an exponential distribution defined as:

$$\phi(\Delta t) = \left(\sum_{i=1}^N \lambda_i \right) e^{-(\sum_{i=1}^N \lambda_i) \Delta t} \quad (4.1)$$

To find Δt , we use the survival function (i.e., the probability that a random variable exceeds a certain value) of $\phi(\Delta t)$, given by:

$$\int_{\Delta t}^{\infty} \phi(t') dt' = e^{-(\sum_{i=1}^N \lambda_i) \Delta t} \quad (4.2)$$

which allows us to calculate Δt as:

$$\Delta t = -\frac{\log(u)}{(\sum_{i=1}^N \lambda_i)} \quad (4.3)$$

where u is a random number drawn from a uniform distribution on the interval $[0, 1]$.

Once Δt is determined, the algorithm identifies which process, denoted as i , has generated the event with a probability given by:

$$\prod_i = \frac{\lambda_i}{(\sum_{i=1}^N \lambda_i)} \quad (4.4)$$

Finally, the algorithm advances time by Δt and repeats the procedure. After an event occurs, it's possible for any of the λ_i values to change.

4.3.2 Generating synthetic timeseries

Once we explained how GilleSpie algorithm works, we present the workflow we used to build our models.

What we want is to implement models who generate synthetic timeseries, similar to the original data. To do this we almost got everything we need, we got the list of the users that interacted on the canvas (not every user interacted), we got probabilities of interaction, we got the number of their interacting pixel and the number of their random position pixel, and so on. But it's only missing one thing to perform this new analysis, that is, the activity rate per hour of every user we are going to involve in these models.

In fact, as explained in [33], we have to simulate N Poissonian processes, that in this case are the placement of the users, everyone with a specific rate of λ_i , ($1 \leq i \leq N$) running in parallel. In our context this rate is the activity rate per hour of every user (basically, the number of placements they made every hour). Once we have all we need we can move on, and explain how our implementation works:

1. calculate the rate per hour of every user;

2. extract random values from a poissonian distribution with mean the sum of the rates of the users we selected. These values are going to be the "timestamps" of the placements of the users;
3. choose a random user based on his rate per hour;
4. assign randomly one timestamp value to the user;
5. now we have 2 possibilities: the user put a pixel in a random position (x, y coordinates), or he put a pixel in one of the 8 surrounding cells of one random user that placed in the last n-minutes (n could be 10 minutes, 1 minute, 30 seconds, or 10 seconds since we studied these timedeltas). This choice depends on the real probabilities of interaction that the user had in the real canvas;
6. Now if the user put a random pixel, he just selects a random color from the 16 color palette, but if it's an interaction pixel then he has to choose between a self interaction or an interaction with someone else. Also this choice is based on probabilities of his type of interactions on the real canvas. Then, choose a random color;
7. In the end, the user and the information related on his placement will be stored in a new dataframe.

We just need to clarify a few things: first of all the possibility of a self interaction could be possible only if the same user is in the list of the users who placed in the last n-minutes.

Second, we have 2 parameters in this situation: the timedelta of the interaction (the "n" parameter) and the length of the list of the users who placed before the iterating user. So, greater these 2 parameters are the most computationally expensive will be the code. For this reason we tried different combinations of theses parameters in the construction of our models.

The last thing we have to mention is that the rate per hour of every user remains the same for different type of definition of interactions, but their probabilities to interact or not change based on the type of interaction we chose.

All of this will lead to different results based on the timedelta we select.

4.4 Partial Information Decomposition

In information theory the sole statistical interconnections that are readily observable through pairwise correlation are limited to bivariate relationships. In the most commonly conducted network analyses, each connection between X_i and X_j is regarded as unconnected to any other connection. There are no straightforward means to deduce statistical associations among three or more variables. Interactions of a "higher order" are established by compiling bivariate linkages in examinations like motifs or

community detection. A significant obstacle to directly investigating higher-order interactions has been the absence of efficient and accessible mathematical instruments for recognizing such interactions.

However, recent advancements in the field of multivariate information theory have facilitated the creation of a multitude of diverse measures and frameworks for capturing statistical dependencies that go beyond the domain of pairwise correlation.

Among the tools that have been applied, one of the most extensively developed is the partial information decomposition (PID) [6]. PID exposes that multiple interacting variables can engage in various distinct information-sharing relationships, encompassing redundant, unique, and synergistic modes.

1. **Redundancy:** Redundant information sharing pertains to data that is "copied" across numerous elements. This implies that the same information could be acquired by observing $X_1 \vee X_2 \vee \dots \vee X_N$, where the variables X_1 through X_N are involved.
2. **Synergy:** Synergistic information sharing involves data that is only accessible when considering the combined states of multiple elements and cannot be obtained from simpler combinations of sources. Synergistic information is only attainable by observing $X_1 \wedge \dots \wedge X_N$, where the variables X_1 through X_N are jointly considered.

Redundant and synergistic information sharing modes can be entangled to form more intricate relationships. For instance, when dealing with three variables X_1, X_2, X_3 , information could be redundantly common to all three, which could be learned by observing $X_1 \vee X_2 \vee X_3$. Additionally, one can contemplate information that is redundantly shared by joint states, such as the information that could be learned by observing $X_1 \vee (X_2 \wedge X_3)$, meaning observing X_1 or the simultaneous state of X_2 and X_3 .

For a finite collection of interacting variables, it is feasible to enumerate all potential information-sharing modes. Given a formal definition of "redundancy," these modes can be computed and quantified.

4.4.1 PID intuition

Williams and Beer [34] introduce a non-negative decomposition of mutual information related to a set of predictor variables $\mathcal{X} = X_1, X_2, \dots, X_n$ with respect to a target variable \mathcal{S} . They break down the total multivariate mutual information, denoted as $I(\mathcal{X}; \mathcal{S})$, into a set of components representing unique, redundant, and synergistic information among all possible subsets of \mathcal{X} .

To achieve this, they examine all subsets of \mathcal{X} , denoted as A_i and referred to as sources. They demonstrate that the redundancy structure of the multivariate information is determined by the "collection of all sets of sources in which no source is

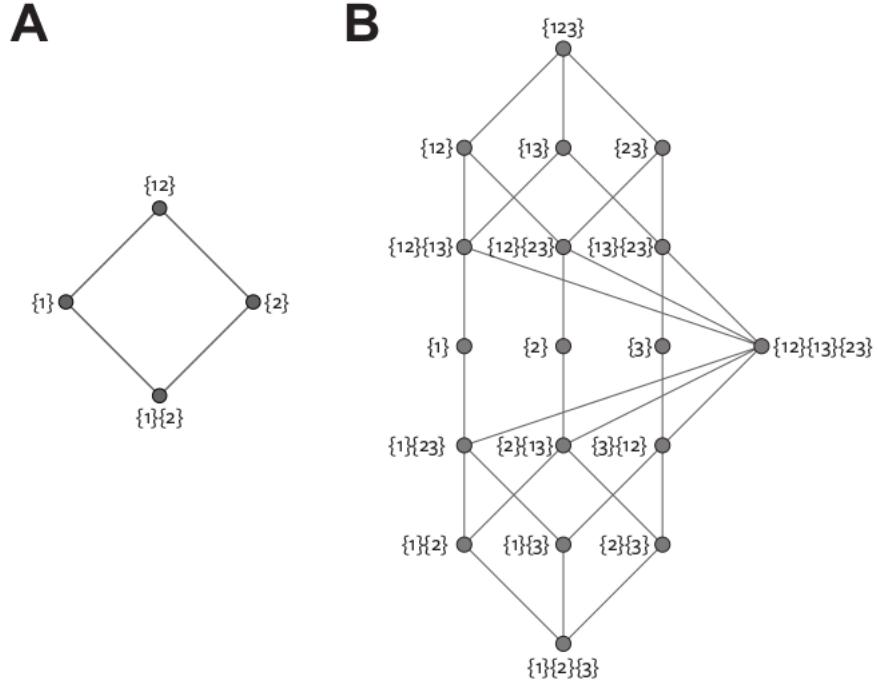


Figure 4.5: Redundancy lattice for (A) two variables, (B) three variables.

a superset of any other" – formally referred to as the set of antichains on the lattice formed from the power set of \mathcal{X} , considering set inclusion, denoted as $A(\mathcal{X})$. In conjunction with a natural ordering, this defines an information redundancy lattice.

Each node within this lattice signifies a partial information atom, and its value is determined by a partial information (PI) function. In Figure (4.5), the lattice structure is depicted for $n = 2, 3$. Each node is described as a set of sources, with sources denoted in braces. The explicit X notation is omitted, so $\{1\}$ corresponds to a source that contains variable X_1 . When a node contains multiple sources, it signifies the redundancy or shared information between those sources. For example, $\{1\} \{2\}$ represents the information shared between X_1 and X_2 .

The information redundancy value of each lattice node, denoted as I_{\cap} , quantifies the total information provided by that node. The partial information function, denoted as I_{∂} , quantifies the unique information contributed solely by that node, whether it's redundant, synergistic, or unique information within or between subsets of variables. For instance, the information redundancy of node $\{12\}$ corresponds to the joint mutual information $I(\mathcal{S}; X_1, X_2)$, and the partial information value of that node represents synergy – information about \mathcal{S} that emerges only when considering X_1 and X_2 together. The partial information value for each node can be computed using a recursive relationship (Möbius inverse) based on the information redundancy values present in the lattice:

$$I_{\partial}(\mathcal{S}; \alpha) = I_{\cap}(\mathcal{S}; \alpha) - \sum_{\beta \leq \alpha} I_{\partial}(\mathcal{S}; \beta) \quad (4.5)$$

In this context, the symbol α , which belongs to the set $\mathcal{A}(X)$, represents a collection of sources. Each source, in turn, is a group of input variables denoted as X_i . These sources are what determine the specific node we're examining.

It's important to highlight that in the context of a two-variable Partial Information Decomposition (PID), there's a connection between the mutual information values for pairs of variables and the terms that describe partial information, and this connection is expressed in the following manner:

$$I(\mathcal{S}; X_1, X_2) = I_\partial(\mathcal{S}; \{1\} \{2\}) + I_\partial(\mathcal{S}; \{1\}) + I_\partial(\mathcal{S}; \{2\}) + I_\partial(\mathcal{S}; \{12\}) \quad (4.6)$$

Consequently, the mutual information between two variables is divided into four components, each of which signifies a portion shared jointly by both variables, a portion exclusive to each individual variable, and a portion that arises from their synergistic interaction. In a similar vein, the mutual information values for each individual variable can be broken down into a combination of shared contributions and unique contributions:

$$I(\mathcal{S}; X_1) = I_\partial(\mathcal{S}; \{1\} \{2\}) + I_\partial(\mathcal{S}; \{1\}) \quad (4.7)$$

$$I(\mathcal{S}; X_2) = I_\partial(\mathcal{S}; \{1\} \{2\}) + I_\partial(\mathcal{S}; \{2\}) \quad (4.8)$$

4.4.2 Partial Entropy Decomposition

A natural extension of the Partial Information Decomposition (PID) is the concept of Partial Entropy Decomposition (PED). Originally introduced by Ince [31], PED diverges from PID in that it doesn't implicate segregating some variables as predictors and others as targets when decomposing the mutual information that a set of predictors transmits about a target.

We use the term H_\cap to represent the total entropy associated with each node within the lattice. This quantifies the complete entropy provided by that specific node. Furthermore, we establish a partial entropy function, denoted as H_∂ , which is structured in a manner similar to I_∂ and is defined through Möbius inversion (as indicated by Equation 4.5). H_∂ serves to gauge the distinct entropy contributed solely by that node, encompassing redundant, synergistic, or unique entropy within various subsets of variables.

In PED, the fundamental approach is quite akin to that of PID. It revolves around the idea that given a redundancy entropy function that adheres to the following axioms:

1. **Simmetry:** $H_\cap(\mathcal{A}_1, \dots, \mathcal{A}_k)$ is symmetric with respect to the \mathcal{A}_i 's.
2. **Self-redundancy:** For single sources, the redundancy is equal to the overall entropy: $H_\cap(\mathcal{A}) = H(\mathcal{A})$.
3. **Monotonicity:** As more sources are added, the redundancy must decrease: $H_\cap(\mathcal{A}_1, \dots, \mathcal{A}_{k-1}, \mathcal{A}_k) \leq H_\cap(\mathcal{A}_1, \dots, \mathcal{A}_{k-1})$

| Node label | Redundancy function | Partial entropy H_∂ | Represented atom |
|------------|--------------------------------|--|--|
| {12} | $H_\cap(\{12\}) = H(X_1, X_2)$ | $H_\cap(\{12\})$ - $H_\cap(\{1\})$ - $H_\cap(\{2\})$ + $H_\cap(\{1\}\{2\})$ | entropy available only from X_1 and X_2 together (synergy) |
| {1} | $H_\cap(\{1\}) = H(X_1)$ | $H_\cap(\{1\})$ - $H_\cap(\{1\}\{2\})$ | unique entropy in X_1 only |
| {2} | $H_\cap(\{2\}) = H(X_2)$ | $H_\cap(\{2\})$ - $H_\cap(\{1\}\{2\})$ | unique entropy in X_2 only |
| {1}{2} | $H_\cap(\{1\}\{2\})$ | $H_\cap(\{1\}\{2\})$ | entropy shared between X_1 and X_2 |

Figure 4.6: Full PED in the case of 2 variables

it's feasible to define a redundancy lattice. This redundancy lattice can then be resolved using Mobius inversion, much like the PID framework. For the case of two variables the nodes of the lattice, their entropy redundancy and their partial entropy values are given in Table (4.6).

So, in direct analogy with the PID we have:

$$H(X_1) = H_\partial(\{1\}\{2\}) + H_\partial(\{1\}) \quad (4.9)$$

$$H(X_2) = H_\partial(\{1\}\{2\}) + H_\partial(\{2\}) \quad (4.10)$$

$$H(X_1, X_2) = H_\partial(\{1\}\{2\}) + H_\partial(\{1\}) + H_\partial(\{2\}) + H_\partial(\{12\}) \quad (4.11)$$

Inserting these partial entropy values into the definition of mutual information we see that:

$$I(X_1, X_2) = H_\partial(\{1\}\{2\}) - H_\partial(\{12\}) \quad (4.12)$$

In essence, mutual information can be understood as the combination of shared or redundant entropy subtracted by the synergistic entropy. In a similar vein, just as interaction information conflates both redundant and synergistic information, mutual information, by its nature, amalgamates both redundant and synergistic entropy components. The exact interpretation of this is difficult, and bases on what "redundant" and "synergistic" entropy is taken to mean.

The partial entropy decomposition shares a similar unique characteristic, namely, it demands the selection of one out of several potential redundant entropy functions. The main possible functions (for a collection of sources (a_1, \dots, a_k)) are:

- *Common surprisal* (by Ince) [35]:

$$H_\cap^{cs}(\mathcal{A}_1, \dots, \mathcal{A}_n) = \sum_{a_1, \dots, a_n} p(a_1, \dots, a_n) h_{cs}(a_1, \dots, a_n) \quad (4.13)$$

- *Shared entropy* (by Varley [6]):

$$h_{\cap}^{sx}(\alpha) = \log \frac{1}{P(a_1 \cup \dots \cup a_k)} \quad (4.14)$$

Common surprisal has been defined as the sum of positive pointwise co-information values, it measures the overall entropy in the node. Basically, local co-information measures the set-theoretic overlap of multiple local entropy values, hence common (or shared) surprisal. Local co-information can be negative, but since the local informations over which the intersection is calculated are positive in this case there is no overlap. So surprisal is local entropy. This redundancy function can return negative partial entropy values, which is difficult to interpret, although Ince et al., argue that a negative partial entropy atom is an instance of misinformation, and can be excluded.

h_{\cap}^{sx} , instead, returns provably non-negative partial entropy atoms, although it is only well-defined for discrete random variables. It measures the amount of information that can be learned about each variable in the node.

Now, once we select a specific entropy function, we move in the same way in both cases calculating the partial entropy value via Möbius inverse:

- Using common surprisal:

$$H_{\partial}(\alpha) = H_{\cap}^{cs}(\alpha) - \sum_{\beta \leq \alpha} H_{\partial}(\beta) \quad (4.15)$$

- Using shared entropy:

$$h_{\partial}(\alpha) = h_{\cap}^{sx}(\alpha) - \sum_{\beta \leq \alpha} h_{\partial}(\beta) \quad (4.16)$$

In this way we calculate all the entropy values for all the atoms in the lattice. For our purposes, as the last thing we need to do, we need to compute the quantity of synergistic and redundant information. To do that, we use Varley definitions [6] of *redundant* and *synergistic structures*:

- **Redundant structure:** when considering higher-order redundancy, we are interested in all of those atoms that duplicate information over three or more individual elements. For example, in the case of 3 elements system, we compute:

$$\mathcal{S}_R = \{1\} \{2\} \{3\} \quad (4.17)$$

Here we used the easier notation in which, for example, $H_{\partial}(\{X_1\}) = \{1\}$.

- **Synergistic structure:** made by all those atoms representing information shared over the joint state of two or more elements. For example, for a three element system:

$$\begin{aligned} \mathcal{S}_S = & \{1\}\{2, 3\} + \{2\}\{1, 3\} + \{3\}\{1, 2\} \\ & + \{1, 2\}\{1, 3\}\{2, 3\} + \{1, 2\}\{1, 3\} + \{2, 3\}\{1, 3\} + \{1, 2\}\{2, 3\} \end{aligned} \quad (4.18)$$

These are the tools we are going to use. However we have to precise that because it's relatively new, the Partial Entropy Decomposition (PED) hasn't been investigated as extensively as the Partial Information Decomposition (PID). While the mathematical foundations are well-established, there are still significant challenges in terms of understanding and interpreting the results. Moreover, applying this framework to analyze complex real-world systems is an active area of ongoing research.

4.4.3 Our PED configurations

To achieve our objectives we decide to use Partial Entropy Decomposition instead of PID, since his application to our case-study seemed to be more appropriate.

Once explained how PED works we can show how this tool has been applied on Reddit Place dataset to find out collaboration patterns.

In section 4.1 we illustrated how to identify a space-temporal window in which an artwork lives in, and thanks to this we are able to work on users that were active in it.

So, we have a list of artwork that we want analyze, what we need to do is understand how to use these information theory measures on it. To apply PED we need a finite set of states for our variables, so we decided to consider only the case of 3 variables systems, because we know that the structure of the lattice mutual becomes even more complex as the context in which it exists expands.

In this approach, our 3 variables system is composed by triplets of users, one for every artwork we consider, who have similar number of placements inside and outside these artworks.

Now, for each of the 3 user in every triplet, we need to give a meaning to his 3 symbols (1,2,3) in manner to calculate the redundant and synergistic information. We tried to use the following different configurations:

- **Spatial cuts & number of times they did anything in intervals of 30 minutes:** basically, knowing the selected artwork's boundaries, we cut it in 2 equal spaced zone, could be up and down, left or right or other 2 different equal zones sliced with 2 different diagonals cut. So, symbols 1 and 2 indicated the number of times every user placed in one of the two zones. Instead, with symbol 3 we collected the number of times users did any placements in a 30 minute time interval for the whole duration of the artwork.
- **Equal or different colors & number of times they did anything in intervals of 30 minutes:** here, symbol 3 is equal to before, but now symbols 1 and 2 indicate the number of pixel with (in the first case), or without (in the second case), one of the same main colors of the artwork they live in.
- **Majority rule & number of times they did anything in intervals of 120 minutes:** for symbol 3 we have a different timedelta (2 hours) and for 1 and 2 we consider respectively the number of "collab" and "conflict" pixel we defined in the previous sections using majority rule.

CHAPTER 4. METHODS

- **Majority rule & number of times they did anything in intervals of 60 minutes:** as before for symbols 1 and 2, but for symbol 3 we reduced the time interval to one hour.

Once, we have defined our symbols configurations, we just have to compute synergistic and redundant information as described by Varley et al. for the chosen configuration. The important thing to understand is the meaning of these information measures which changes depending on the configuration chosen and the context in which the measures are applied. We will discuss about this in the results section.

Chapter 5

Results

5.1 Exploratory data analysis

As we said before the dataset is composed by the following quantities: **time**, **user**, **x-coordinate**, **y-coordinate** and **color** and contains over 16M placements.

Our first aim was to find the more relevant quantities that could make us able to find collaboration patterns, so we started to dive into the data to discover something useful for our purpose.

5.1.1 Activity trend analysis

In this scenario *activity* tells us a lot about users behaviour; the term stands for the presence of the users on the canvas, obviously in terms of their placements, of the coordinates of their placements, of the chosen colors of them and so on.

As a preliminary step, we examined the progression of activity throughout the timeline. This assessment encompassed the quantity of placements, the involvement of active users, and the evolving use of the canvas over time. It is crucial to recognize that activity levels are not uniform across time intervals. Therefore, it is intriguing to trace the pattern of activity as the canvas, initially void, gradually transforms into a complete image.

As this progression transpires, users inevitably find themselves forced to engage with one another due to the increasing scarcity of empty canvas spaces. To conduct this analysis, we segmented the dataset into discrete time segments spanning an hour each. Within these segments, we tabulated the counts of interactions, users, and non-blank pixels. These counts were then graphically represented over the timeline to reveal the dynamic evolution of activity (5.3).

Also, there exists an opportunity to generate a visualization that captures the activity distribution and patterns across the canvas throughout the entire event. This visualization takes the form of an activity heatmap, which highlights the spatial distribution of user interactions. This is achieved by constructing a two-dimensional histogram counting the activity level at each individual pixel.

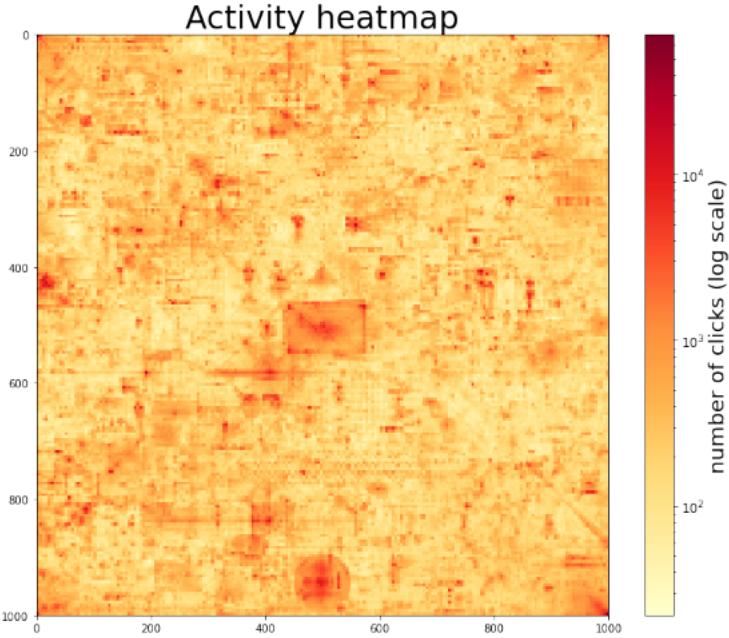


Figure 5.1: Activity heatmap of the entire experiment aggregated on a grid of 5-pixel per side for visualization.

For illustrative purposes, a concrete demonstration of this concept is presented in (5.1). The depicted heatmap showcases the extent of activity throughout the entirety of the experiment. In order to facilitate visualization, the values denoting activity are aggregated in a spatially simplified representation of the canvas. This representation is achieved by dividing the canvas into square windows, each encompassing a group of 5 pixels on each side. This spatial coarsening aids in revealing the overarching trends of activity distribution across the canvas.

Our dataset contains more than 1 million unique users who, at least, placed 1 pixel during the whole experiment. Obviously the minimum number of placements of a single user is 1, the maximum number is 656, and the mean number is about 14 placements.

But, to have a more accurate view about the number of placements of every user we can have a look at (5.2): we can see that the majority of user just placed no more than 5 pixel. Maybe because people were curious about this social experiment but didn't have a lot of time to spend on the canvas since you have to wait 5 minute to place the following pixel. Then we can see a solid population of users that placed between 10 and 100 pixel during the 3 days. Instead, the number of people persistent on the canvas, capable to put more than 100 pixel, is lower than the 2 previous "categories" but still consolidated. Probably these were the committed users.

In Figure (5.3 (a)), we can observe a pattern in activity levels over time during the event. The number of clicks and active users per hour undergoes periodic fluctuations, roughly following a 24-hour cycle that aligns with the circadian rhythm of the American continent. Since a significant portion of Reddit users are from the United States, this

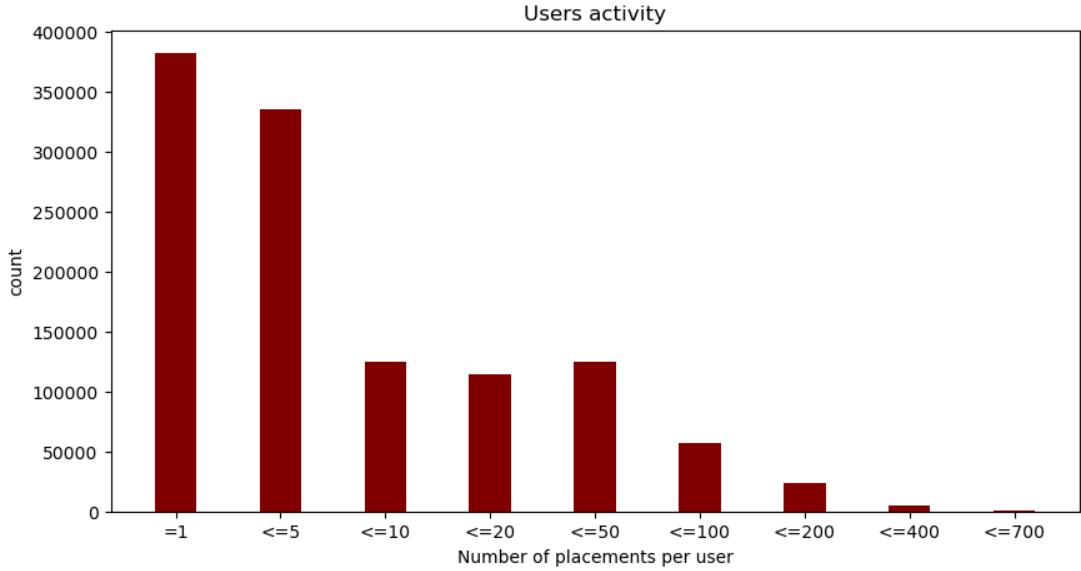


Figure 5.2: Number of placements per user distribution

cycle makes sense. Additionally, there is an uptick in activity after the first day, which can be attributed to the spread of information about the event.

Turning our attention to the canvas's evolution over time (as depicted in (5.3(b))), we notice a notable trend. After approximately the first 30 hours, the amount of empty space on the canvas significantly diminishes, with most of the available space being utilized. The most substantial changes on the canvas occur during the initial day when users fill in the empty spaces. Subsequently, due to the limited available space, users are compelled to engage more with each other within existing artworks and in neighboring areas. This phenomenon is clearly visible studying the timelapse of the experiment. The most significant canvas transformations happen within the first day and a few hours. After that point, due to space constraints, pre-existing structures either undergo remodeling or become more consolidated.

Now, thanks to this result and after that the number of placement per user distribution has been studied we can say that each user faces a practical limitation on their activity, primarily imposed by interaction frequency restrictions, typically ranging from 5 to 20 minutes. Consequently, the majority of users perform a limited number of updates, with the mean value being around 14 updates per user and a maximum on 656.

5.1.2 Inter-event times analysis

Analyzing inter-event times is an indispensable tool for gaining a deeper understanding of how users interact with systems and content, making it a fundamental component of data-driven decision-making processes.

We started this section by studying inter event times distribution for the 10 most active users on the canvas (5.4). We selected these 10 users with a total of 5778 pixel,

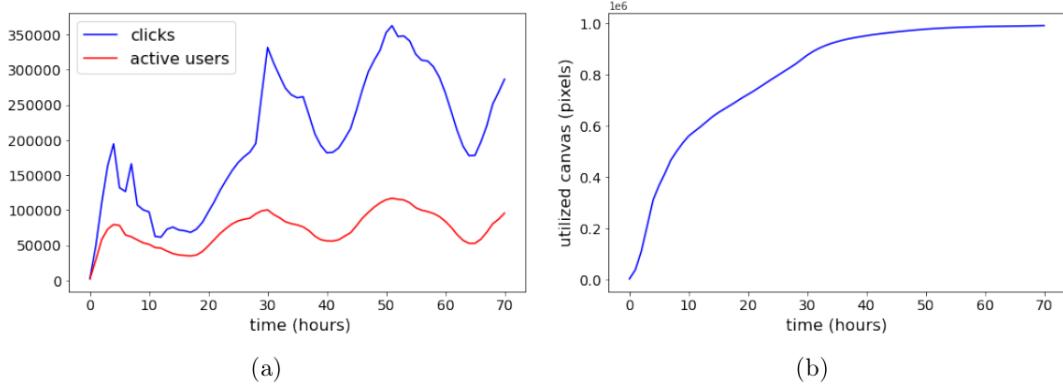


Figure 5.3: (a) Activity over time: clicks per hour (blue) and active users per hour (red). (b) Canvas coverage over time.

and mean number of placements about 577.

The peak of the distribution is about 20 seconds between a pixel and the following; this is an interesting result considering that this analysis involves all the days of the experiment, indicating that probably these users were in someway coordinated. Also, the distribution seem not to be Poissonian.

In addition, we compared the PDF of inter event times user per user in the whole canvas and the PDF of inter event times selecting only neighbours users of the iterating pixel in a timedelta of 10 minutes (5.5), basically in the latter we analyzed interacting pixel.

The plot on the left confirms how people must wait 5 minutes to place the next pixel, but also we can see that there is a minor population who didn't respected this 5 minutes bond, probably because they were "super-users" or because in that time there were an inferior timedelta than the 5 minutes one. Instead, the plot on the right is very useful in terms of interactions, in fact it says that, in a timedelta of 10 minutes, people mostly interacted with the first 60 seconds of the time interval.

This last result brought us to test different timedeltas for our concept of interaction.

In addition, this kind of analysis was also useful in comparing what happens inside and outside the artwork. In fact, we selected a list of 16 different artworks on which to perform different types of analysis, and for each of these we selected a triplet of users who had similar activity quantitatively both inside and outside.

In this case, we computed the inter event time distribution for triplets of users (5.6). The graph therefore shows that the activity of the users inside the artworks is tighter than outside the artwork, indicating that they are more interested in creating something meaningful only in that space-time region.

5.1.3 Activity through the time

This analysis helps in identifying patterns, detecting anomalies, and discovering differences between users behaviour during the whole event.

CHAPTER 5. RESULTS

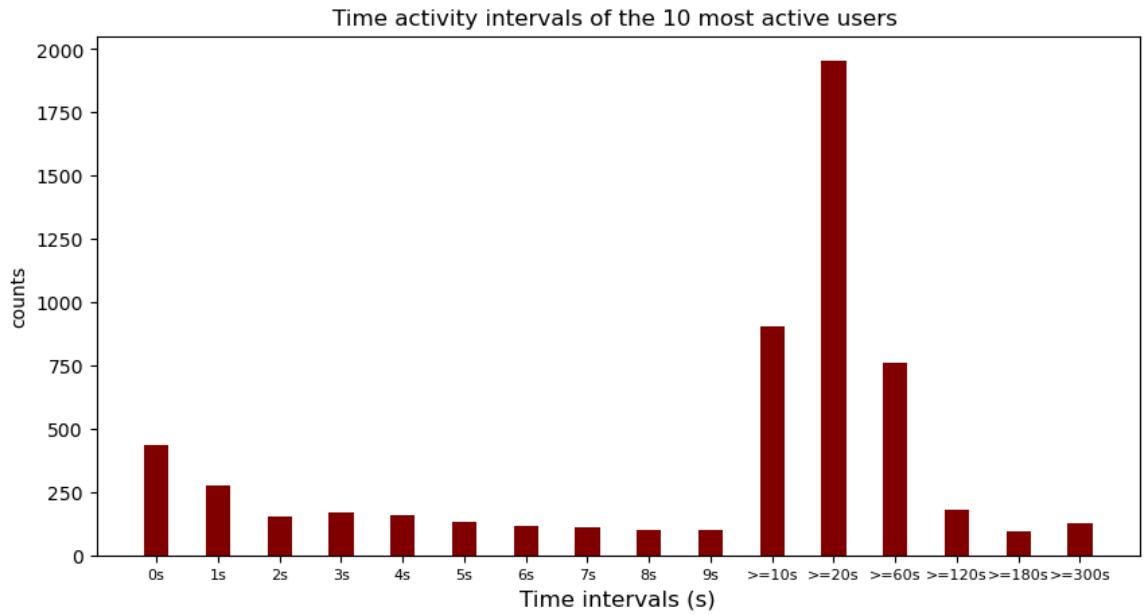


Figure 5.4: Inter placement times of the 10 most active users on the canvas.

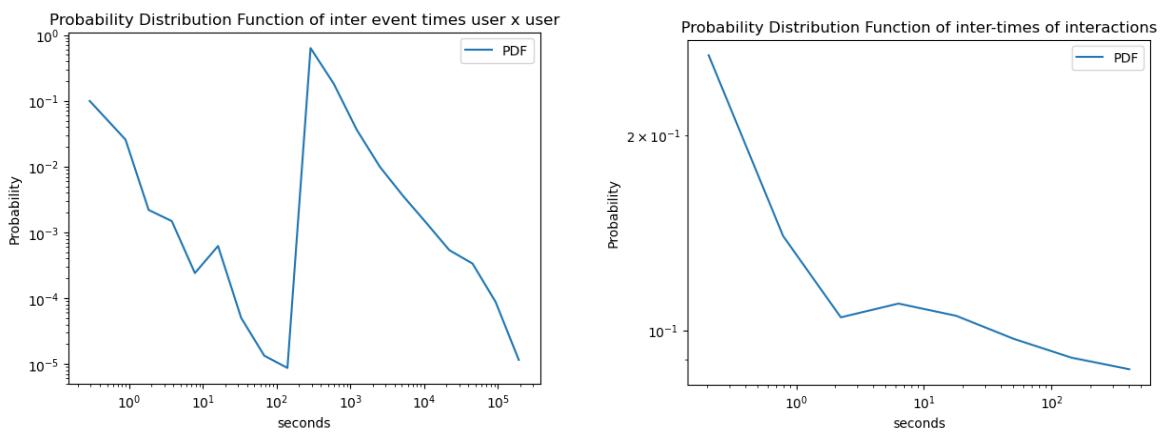


Figure 5.5: On the left: PDF of inter event times of every user on the canvas. On the right: PDF of inter event times of only neighboring pixel, i.e. interacting pixel.

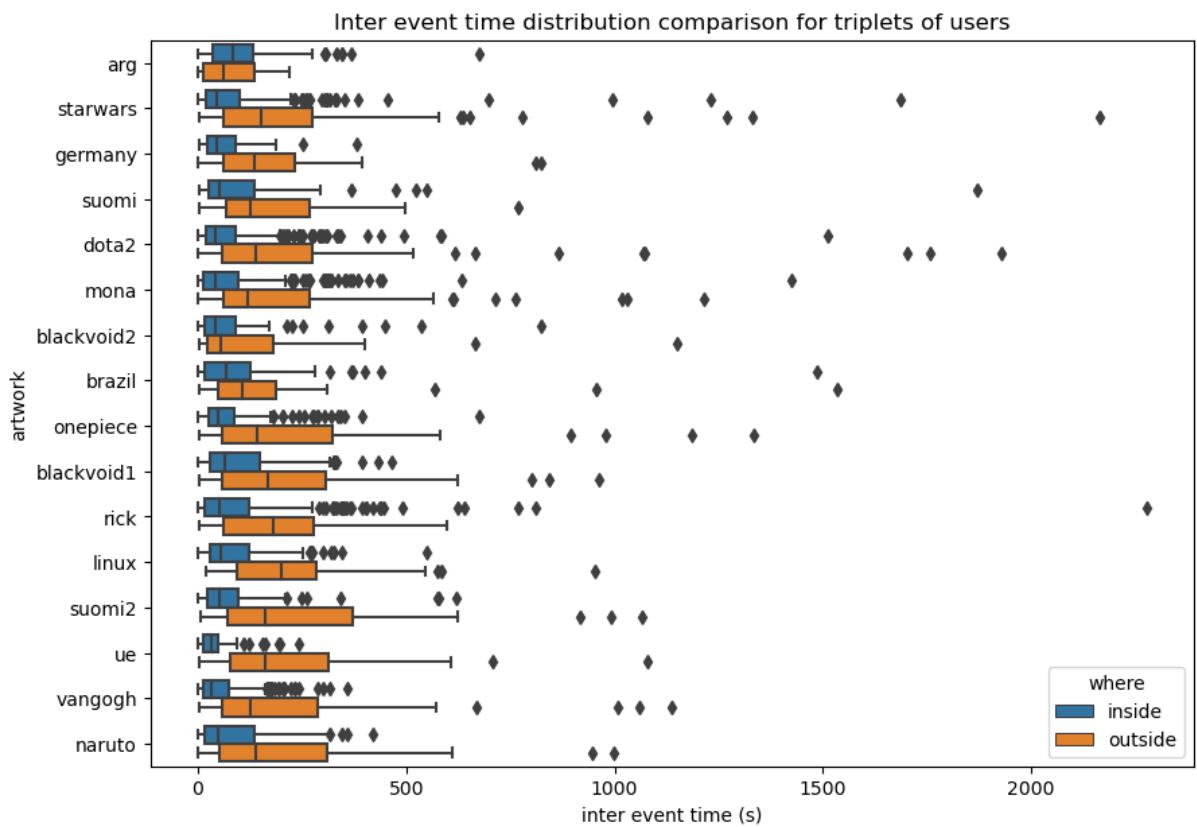


Figure 5.6: Boxplot distributions of inter-event times inside and outside some artworks considering the activity of one triplet of users for each one.

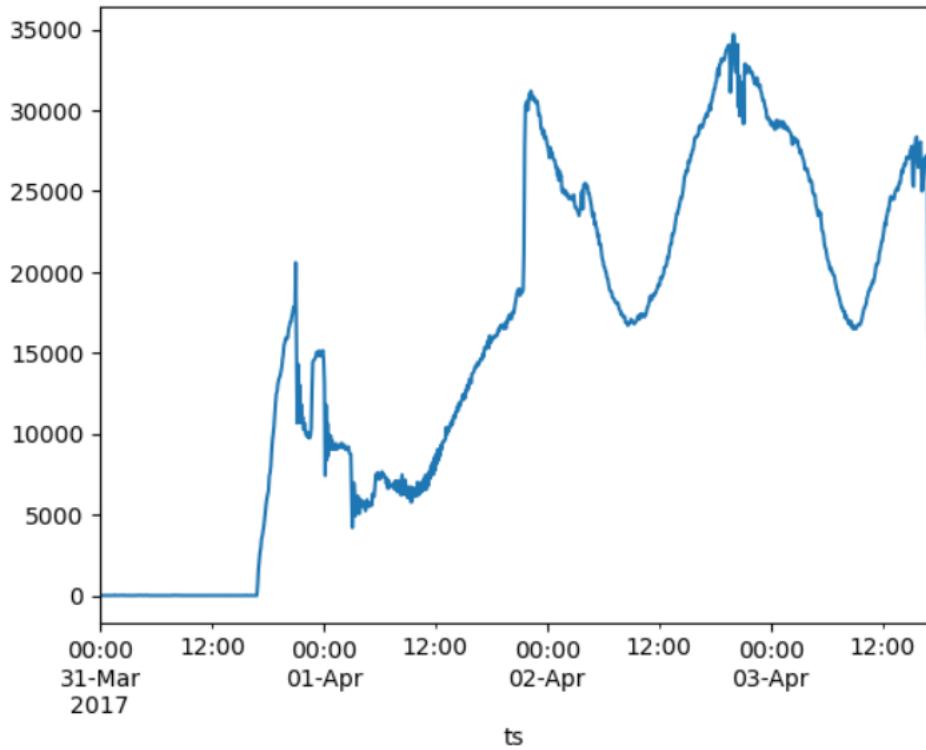


Figure 5.7: Activity trend in the whole experiment divided in 5min intervals.

So, after inter-event times we also analyzed the activity trend through the time first in the whole experiment (5.7), and then day by day (5.8).

In the first case the activity is updated by 5 minutes intervals. The graph shows how the activity is increased during the 3 days period. In the second case we have a specific evaluation of the activity during the single day separately. At the beginning just a few people knew the existence of the social experiment, instead on the last day users seem to be very present on the canvas, probably the committed ones. The graph (5.7) shows how at the beginning we had just a few users active on the canvas, probably because not a lot of people knew about the experiment. However since the 1st of April, we can see how the experiment gained media relevance showing a concrete increase of the activity.

Plus, observing to (5.8), besides the first day, we can see that the last 3 days of the experiment have a similar trend especially the last 2, and also there is a specific time interval (between 6 and 12 AM) in which the activity decreases respect to the previous and the following hours.

5.1.4 Color feature analysis

Finally, the "color" feature of the dataset is one of the most important to understand the meaning of a single pixel: if 2 different colored pixel are neighbors maybe the users

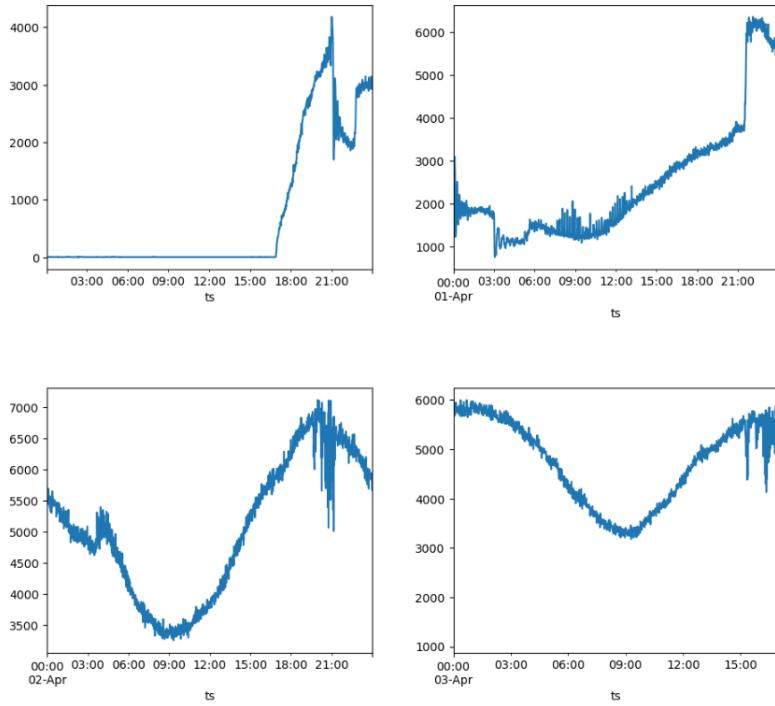


Figure 5.8: Activity trend day by day divided in 1min intervals.

who placed them are not collaborating, or maybe yes, or maybe the fact that they are neighbors is just a consequence of the huge quantity of users on a limited space.

So, we don't know the meaning immediately, probably we need a deep study of the zone or something else, but know their colors facilitate for sure our analysis. That said, color allows us to understand something more about the meaning of a single pixel (or a conglomerate of them) placed by users.

So, for this reason we studied which colors are the most frequent to see on the canvas (5.9). The plot shows that black and white are the most chosen colors, probably because they have more applications than a pink color or an orange color. Also blue and red seem to be very used.

5.2 Identifying collaborations and conflicts

Once we have explained in methods section how to identify and manage artworks we can perform different things on them, for example we can try to label users as collaborators or not on a specific artwork.

Thanks to [28] we understood how people organized their actions on the canvas by using threads on their subreddit, but also this article is useful because, through the study of different topics of specific subreddit, it identifies peaceful (and not) relations between several neighbors artworks.

In detail we discovered a specific alliance (5.10) between "argentina", "brazil", and "suomi", whose also were in peace with "SouthAfrica" and "SquareSpiral".

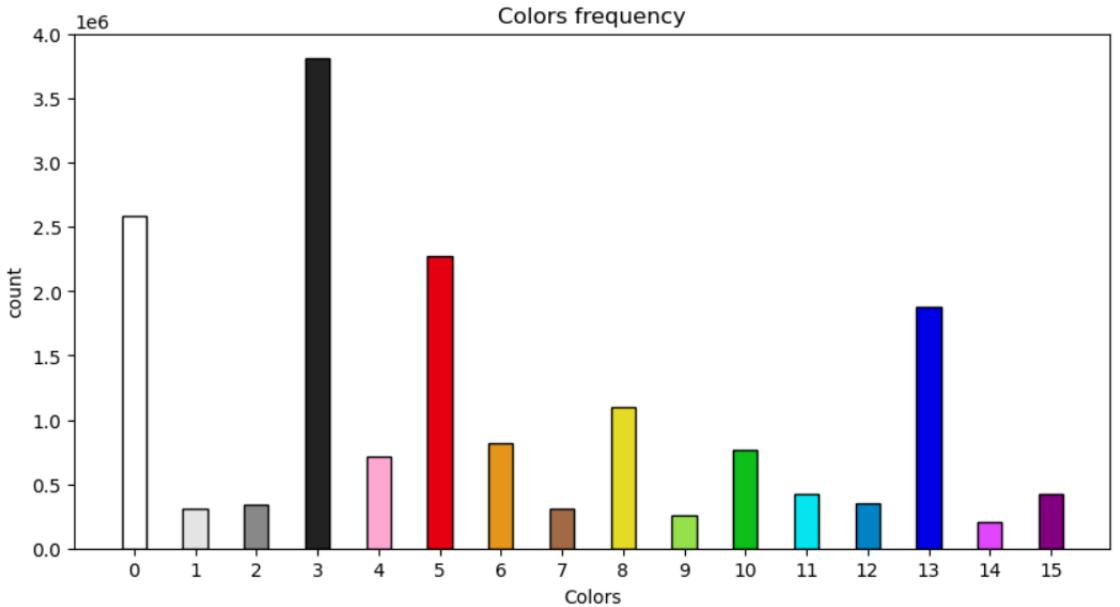


Figure 5.9: Colors distribution.

For our purposes we are going to analyze just the collaboration between "argentina" and "suomi" for a specific reason: although all of the artworks are in neighboring locations, a more efficient analysis can only occur for those 2 as they are at the same y-coordinates thus making them perfectly neighboring.

Also thanks to the timelapse of the experiment we discovered the presence of a specific conflict on the canvas between "germany" and "france" (5.11).

The results of this analysis are stored in (5.15). The term "collab" in this case refers to the users who only placed pixel with the main colors of their flag, for example in the case of Suomi flag we have white and blue.

We start by describing what happened between Suomi (a) and Argentina (b): as previously stated, these 2 artworks were neighbors (one next to each other) and in a peaceful relationship, in fact what we see is that in both cases the neighbor population behave in a respectful manner on their flag, positioning mainly pixel with the colors of the flag in which they are.

This means that they helped their neighbours to maintain the definition of the flag on the canvas during the whole experiment despite the fact that it was not their flag.

Different is the behaviour we see studying what happened between Germany and France (c): we know that they were in a conflict that France lost and for this reason it was forced to move in another area of the canvas. For this reason we expect to see something different from the two other cases, in fact the French users, who were committed to build their own flag and were active also on the German flag, placed just 521 pixel with the colors of the German flag and more than 2500 pixel with different colors of it.

CHAPTER 5. RESULTS

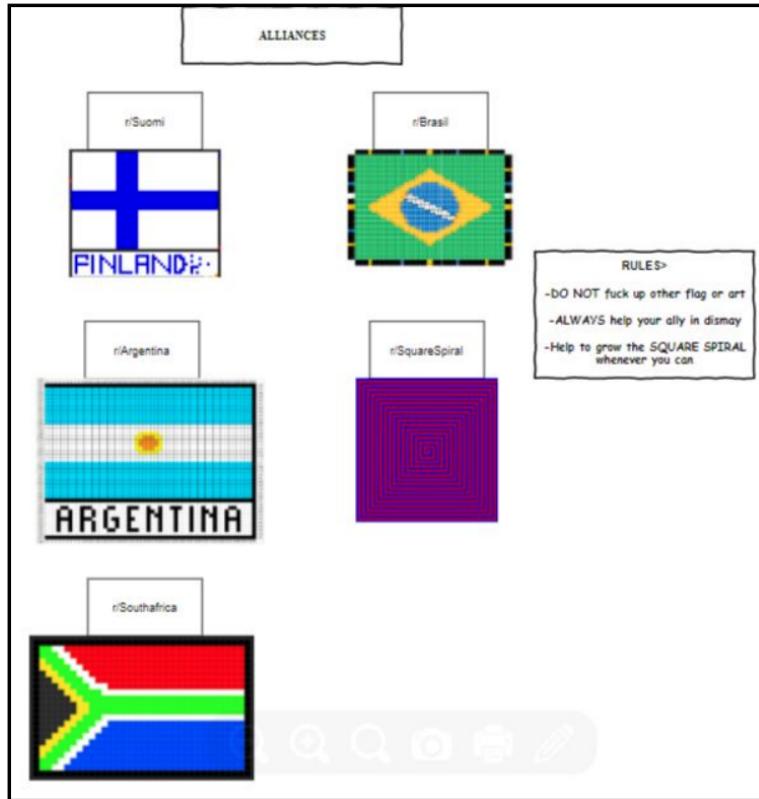


Figure 5.10: Peaceful relationships scheme between specific artworks

This result could be interpreted in the following way: of course french users wanted to destroy German flag, in fact they seemed to be very committed in that placing more than 2 thousand pixel with different colors from the flag, but why we observe a consistent number of pixel with the same colors of the flag (521)? Well, this could be explained considering the fact the 2 flags have a color in common, that is red, so when we see that number we have also to think about the fact that in the same area the same french users tried to safe their flag without success.

This analysis was useful to show ho people in different areas on the canvas has behaved based on the relationships they have with the communities that they had as neighbours.

That said, we now know that there has been different behaviours in different areas of the canvas, someone has respected the neighbour's space boundaries, someone less.

For this reason we can use these information to try use a more powerful tool to label users as "collaborators", as them who just placed pixel to reinforce their baseline artwork, and "attackers", as them who just placed pixel that are not helpful in defining the image of the artwork.

To do this we used a **majority rule**:

- We selected an artwork and the users that were active on it;



Figure 5.11: The conflict under our analysis: germany forced france to move from their original position.

- We removed from our analysis the users who had less than 3 placements, because we are looking for committed people;
- Then we label the users using a majority rule (5.16): if the color of the central cell is equal to the mode of the neighbors colors then that pixel is a "collab" pixel, otherwise is "conflict" pixel;
- Finally we can label users: users with 0 "collab" pixel are **attackers**, users with 0 "conflict" pixel are **collaborators**.

The majority rule we explained here is a powerful tool: the main problem with the "color" feature of the dataset is that we can't know if that color is useful for the artwork in which it was situated or not. Just think to a flag artwork, for example germany flag, let's say that we are studying a yellow pixel, one might be led to think that it is a flag helper pixel, but perhaps it is actually in the area where the red color of the flag would normally go. This is just an example; other scenarios may occur, but it gives the idea.

So, to perform an analysis using this rule, we take back the two situation we studied before, Argentina and Suomi respecting each other, and France and Germany making war on each other. Also, we have to say that, the area that housed the Germany artwork after a few hours will go to host a new artwork, that of the EU, probably indicating the presence of a new conflict, hence we studied this new transition too.

Results are stored in 5.1. We collected the duration of these artworks, the number of total users in them, and also we labeled the people based on the majority rule we mentioned before selecting only the users with more than 2 placements on the baseline artwork.

Again, Argentina and Suomi had similar situations: same duration of 25 hours, similar number of users, but the ones of Argentina seemed to be more active than the ones on Suomi looking at the column "more than 2 pixel users". What we are interested more are attackers and collaborators, and we can see as in both flags the collaborators are the majority of the population, indicating a greater tendency to build the artwork with the aim to maintain his definition till the end.

CHAPTER 5. RESULTS

| Suomi Flag | |
|--|------|
| Total users | 5574 |
| Arg "collab" on the flag | 361 |
| Pixel placed with the same colors of the flag | 1063 |
| Pixel placed with different colors of the flag | 28 |

Figure 5.12: (a)

| Arg Flag | |
|--|------|
| Total users | 5802 |
| Suomi "collab" on the flag | 324 |
| Pixel placed with the same colors of the flag | 1465 |
| Pixel placed with different colors of the flag | 63 |

Figure 5.13: (b)

| Germany Flag | |
|--|-------|
| Total users | 23315 |
| French "collab" on the flag | 394 |
| Pixel placed with the same colors of the flag | 521 |
| Pixel placed with different colors of the flag | 2649 |

Figure 5.14: (c)

Figure 5.15: Users behaviour analysis on different types of artworks. Suomi and Argentina were neighbors on the canvas in a peaceful relationship and this is verified by their users behaviour (a), (b). France and Germany were in a conflict, as we can see by data in (c).

Different is the behaviour of France artwork, who just lived 4 hours, but with more active users on it than Argentina and Suomi. Anyway the committed population of these users is about 2307, some of these were attackers, others (the majority) were collaborators. The 4 hours selected for the duration of artwork France indicates precisely the life of the artwork until just before it was defeated, consequently the interesting figure here is that of the number of attackers, which, although less than the number of collaborators, was able to change the flag of rivals at will. Probably the Germans in this case were better organized and perhaps even encountered a France that did not expect such a conflict, almost to the point of seeming unprepared to face them.

The last transition we have to comment is from Germany to EU: German artwork lasted for 20 hours, and in it were more than 20 thousand users but only 6717 were the most active on it. Interesting is that the number of attackers is almost the same of the one of collaborators, anyway we know that Germany left his area to the European Union but we don't know how. The indication about the number of attackers and collaborators is telling something different from the France-Germany conflict: probably could be that at the beginning Germans struggled to maintain their artwork but then agreed with other users to leave the space building a flag involving more than a country with the aim to be more powerful than being just them. This is just a guess, but could

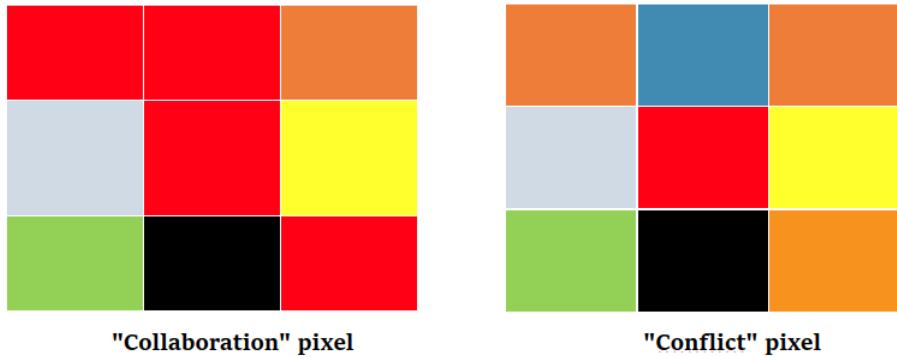


Figure 5.16: Majority rule examples: on the left an example of collab pixel, on the right an example of a conflict one.

| | Duration | Total users | More than 2 pixel users | Attackers | Collaborators |
|--------------------------|----------|-------------|-------------------------|-----------|---------------|
| <i>Argentina</i> | 25 h | 5802 | 1756 | 276 | 479 |
| <i>Suomi</i> | 25 h | 5574 | 908 | 101 | 538 |
| <i>France to Germany</i> | 4 h | 9006 | 2307 | 637 | 992 |
| <i>Germany to EU</i> | 20 h | 23315 | 6717 | 1423 | 1437 |

Table 5.1: Users analysis in different types of artworks: argentina and suomi reached an agreement of peace, instead of France who was involved in a conflict with Germany. At the end also Germany lost his position, transforming itself into EU artwork.

explain the numbers we encountered.

However, it should be pointed out that we don't know anything about the users that were neither collaborators nor attackers, in fact in some cases, these probably played a decisive role within the respective artwork.

This rule, therefore, does not give us the full picture of the situation, but it certainly helps to understand what is happening on a local scale

While it is true that even our definition of majority rule may not be error-free, we must say, however, that it is a much more stringent criterion than just color analysis, as one might be led to do what he sees around him.

5.3 Interactions analysis

We mentioned, in the previous chapter, that our definitions of interaction are all based in terms of proximity in time and in space, to achieve a solid representation of what is around a specific placement.

Anyway, in here, when we find an interaction we don't know yet if it indicates

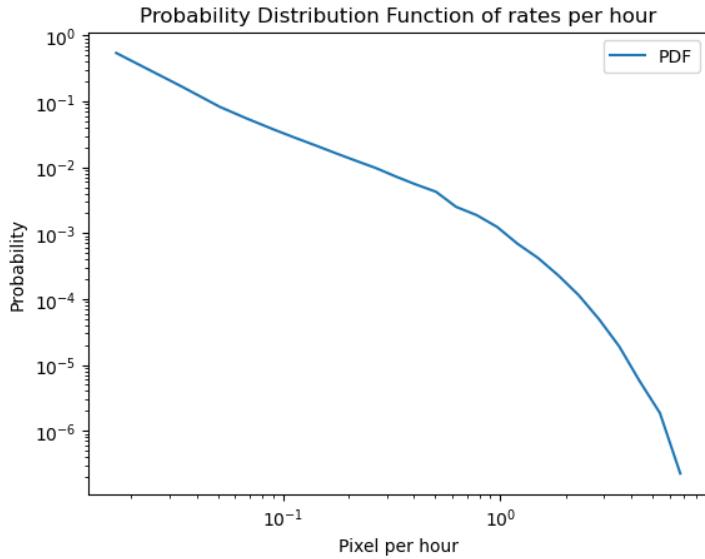


Figure 5.17: PDF of rates per hour of every user

a collaboration or a conflict between 2 users or more, because we don't have any indication about the colors, but for sure we can conclude that something between has happened since they have been in contact.

This encounter between users could be done on purpose, that is with the specific intent to put a pixel next to someone else, or randomly, because of the limited space in which is possible to place a pixel. Our objective is to distinguish between random interactions and intentional ones while identifying user tendencies towards interaction propensity.

Also, as explained in the previous chapter, since we don't know when, in a window of 10 minutes, the users are more inclined to interact, we tried different timedeltas for our concept of interactions to discover the main differences between our definitions.

Before of discussing our findings, let's point out something more: in figure (5.17) we can have a look at the PDF of the rate per hour of every user on the canvas. This is a long-tail function, his shape is caused by the fact that there is a limited number of pixel that every user can place in an hour because of the presence of a 5 minutes bond. Now, knowing how the experiment has lasted and how many pixel are present on the canvas, we computed the number of pixel that are placed every hour on the whole canvas, we are talking about 28000 placements per hour. This is almost the 3% of the size of the entire canvas.

Now, we know that these placements could generate an interaction or not, so knowing the probabilities of the users to interact or to put a random pixel and considering a timedelta of 10 minutes for our interaction (first definition), we found out that every 10 minutes we have the $\frac{2}{3}$ of pixel contributing to interactions and $\frac{1}{3}$ placed randomly. This means that every 10 minutes we have a about 18216 interacting pixel.

Of these 18216 pixel, not all of them cause intentional interactions; some may result from the canvas size or other random factors. To determine how many interactions are

CHAPTER 5. RESULTS

intentional, we use a null model, reshuffling user positions randomly and calculating interactions. We then subtract this from the actual canvas interactions, revealing the 'real number of interactions' made intentionally.

Next, we divide this number by the total number of pixel placements to obtain interactions per pixel. By multiplying this rate by the number of pixels created every 10 minutes, we find that there are approximately 150,000 interactions made every 10 minutes. This suggests users are more inclined to interact intentionally.

So this is what's happened considering the case of 10 minutes, but to have a complete picture of the situation let's have a look to table (5.2).

| | 10 m | 1 m | 30 s | 10 s |
|--|---------------|-------------|-------------|-------------|
| <i>Number of interactions</i> | 132'689'001 | 15'446'996 | 8'168'274 | 2'989'010 |
| <i>Random interactions</i> | 9'727'104(7%) | 985'622(6%) | 501'898(6%) | 177'992(6%) |
| <i>"Real" interactions</i> | 122'961'897 | 14'461'374 | 7'666'376 | 2'811'018 |
| <i>Mean probability per user to put a pixel to interact</i> | 66% | 25% | 17% | 8% |
| <i>Mean probability per user to put a random pixel</i> | 34% | 75% | 83% | 92% |
| <i>Mean number of pixel per user to generate at least an interaction</i> | 9.13 | 3.57 | 2.39 | 1.11 |
| <i>Mean number of random pixel per user</i> | 5.01 | 10.57 | 11.75 | 13.03 |

Table 5.2: Interactions analysis for different timedeltas on the original canvas. Random interactions refer to the number of interactions we found using the respective null model.

First of all, considering the 10 minutes timedelta, we see a huge quantity of interactions, about 130 millions. Of these, the 7% are probably caused randomly, obtaining anyway a big quantity of real interactions. In this case, users seem to be more prone to interact looking at their probabilities. On mean, we have that every user got 9 pixel that generate at least an interaction against 5 pixel who are placed randomly.

The timedelta of 10 minutes is the only case in which users tendency to interact is greater than their tendency to act randomly. In fact, moving to inferior timedeltas we see how the probabilities of interaction are always lower than the random position ones. We practically observe a reversal of the users activity trend.

As a manifestation of this we can take a look at figure (5.18), representing the PDFs of pixel per user who generated at least an interaction for different timedeltas; reducing the time-interval, distributions stop earlier than the previous, indicating that the number of interacting pixel is lower.

Also is interesting to see that the number of users who placed pixel to interact decreases with the decrease of the timedelta, in fact for 10 minutes timedelta we have

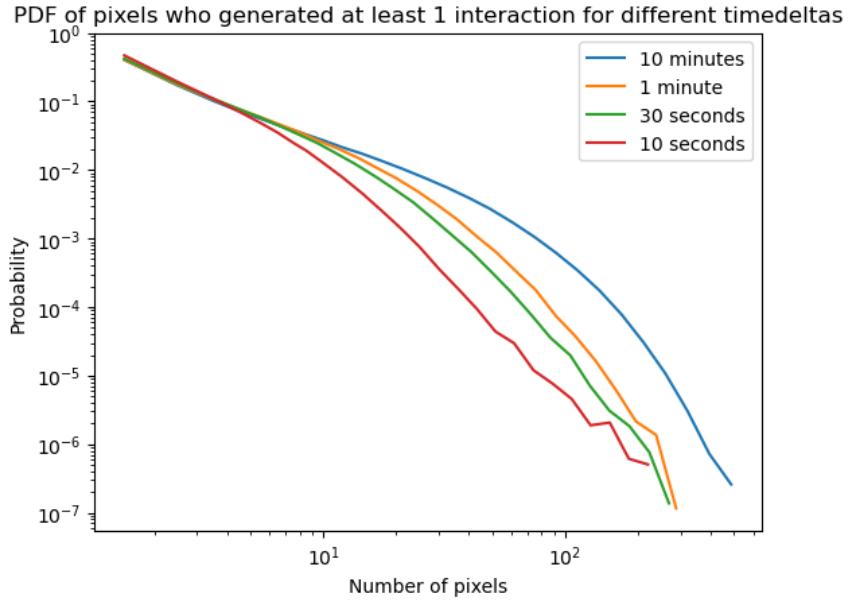


Figure 5.18: Probability distribution functions of pixel per user who generated at least an interaction

782543 users, for 1min we have 576053 users, for 30sec timedelta we have 497182 users, and for 10sec timedelta we just have 361580 interacting users.

An interesting thing to notice is that the rate of random interactions we discovered by the null models is constant at 6-7% of the total interactions through the timedeltas. This indicates we probably have correctly circumscribed the random effect.

Another curious finding is that, looking at the number of total interactions, if we reduce the interval from 10m to 1m, there are not 10 times fewer, or from 1m to 10s, there are not 6 times fewer. This means that probably these events are concentrated in the last moments of the time interval, i.e. we observe a phenomenon that is not linear in the timedelta, probably in the last moments of that specific interval there are more.

After this findings become natural to ask ourselves if the most active users are effectively the ones who interact the most. This is explained in figure (5.19). We selected the 100 most active users on the canvas, we collected their number of interactions for different timedeltas and we made a scatter plot of their correlation. The only case in which seem to be a low positive correlation between the two features is in the 30 seconds timedelta interaction. Probably we identified the most likely time interval of interaction, since one is led to think that the most i am active on the canvas and the most i will interact with someone else both because i want to do that and because of the presence the random effect.

5.4 PED analysis

In this section we show the results of PED applied on artworks and users using different configurations. We started by working on real data, we compared what happen com-

CHAPTER 5. RESULTS

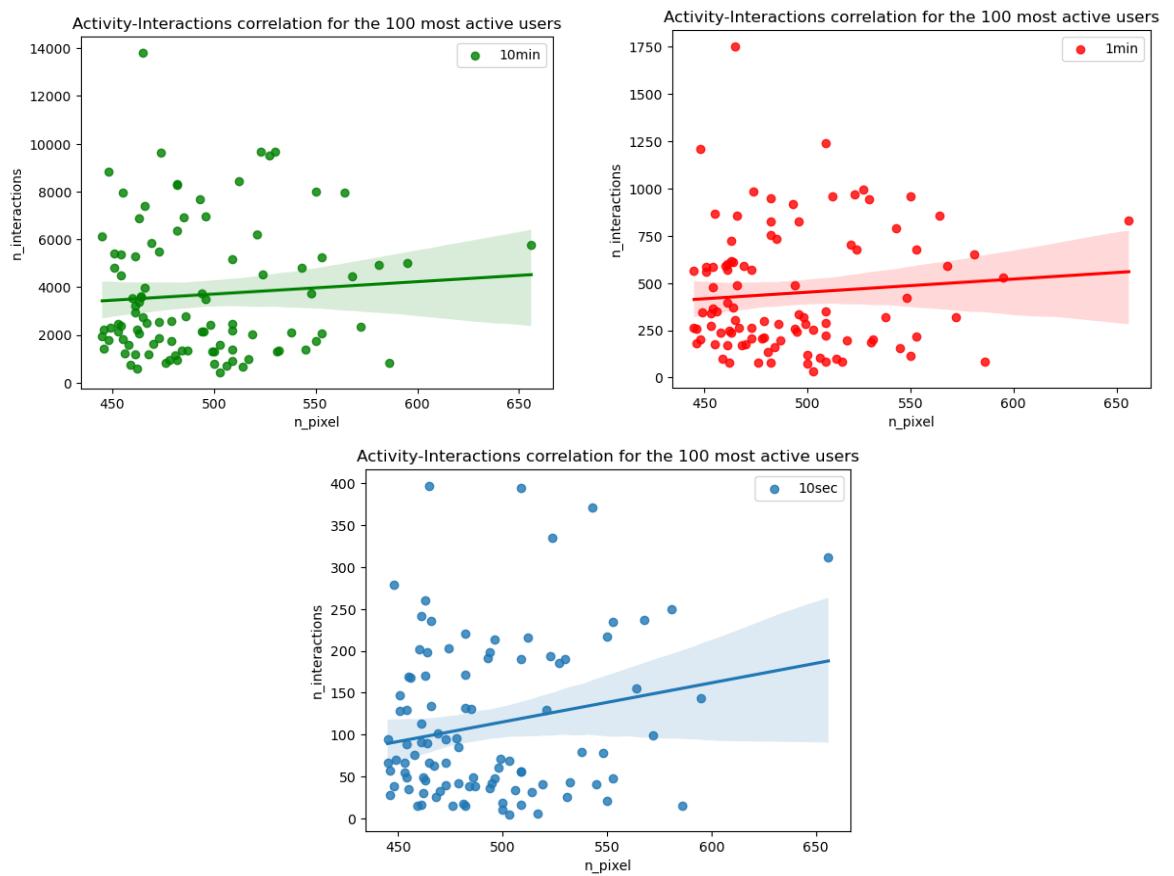


Figure 5.19: Activity-interactions correlation for different timedeltas

CHAPTER 5. RESULTS

paring users activity inside and outside the artworks, then we created some synthetic timeseries to confront PED results between them and real data, and finally we tried different definitions for the symbols of PED to find out which is the best configuration for our purposes.

Before showing our results, it is necessary to clarify what synergy and redundancy symbolize in our case study. Redundant information refers to the part of the information that can be predicted or inferred from other variables in the system. In other words, it is the information that does not provide new insights or surprise once other relevant variables are known. Synergistic information, on the other hand, refers to the information that arises from the interaction or combined effect of multiple variables in the system. It represents the additional information that cannot be predicted or attributed to any single variable alone. Synergy implies that the collective behavior of the variables leads to emergent patterns or effects that go beyond what each variable contributes individually.

In our situation redundancy is what we expect to be higher inside the artworks because it represents the fact that users are working on the same symbol (that is a specific action) and we know that in some way users collaborated inside these zones. Synergism, instead, stands for the fact that users, in parallel, work on different commitments, maybe different areas, or put different colors among them, and so on; so we don't really know what to expect in terms of synergism, it is probably possible that outside the artworks this measure is higher than inside because they couldn't have the same objective as inside of it.

That said, we are ready to explore our results.

5.4.1 Inside & outside the artworks

To initiate our analysis, we initially examined variations in information measures within the context of artworks and compared them to those outside of artworks. Our study involved a selection of 18 distinct users, resulting in the generation of 10 unique user triplets. In some cases we studied just one triplet per selected artwork, in others more than one triplet of users per artwork. Notably, this analysis encompassed five specific artworks, namely Linux, Suomi, Brazil, Mona, and UE.

For our computations of information measures, we adopted the definition of PED ("Spatial cuts & number of times they did anything in intervals of 30 minutes") as detailed in Section 4.4.3.

This approach brought us to have 80 different information values (40 synergistic and 40 of redundancy). Figures (5.20) and (5.21) show our findings. Synergism seems to reach higher values than redundancy either inside and outside the artworks.

Also, the different typology of cut doesn't affect the measures, in fact the trend for the specific triplet is similar for all the cases, so maybe this definition of PED is not the best way to map users spatial areas of placement.

Finally, probably the most important result is that measures outside the artworks are higher than inside them, and this is interesting especially for the case of redundancy.

CHAPTER 5. RESULTS

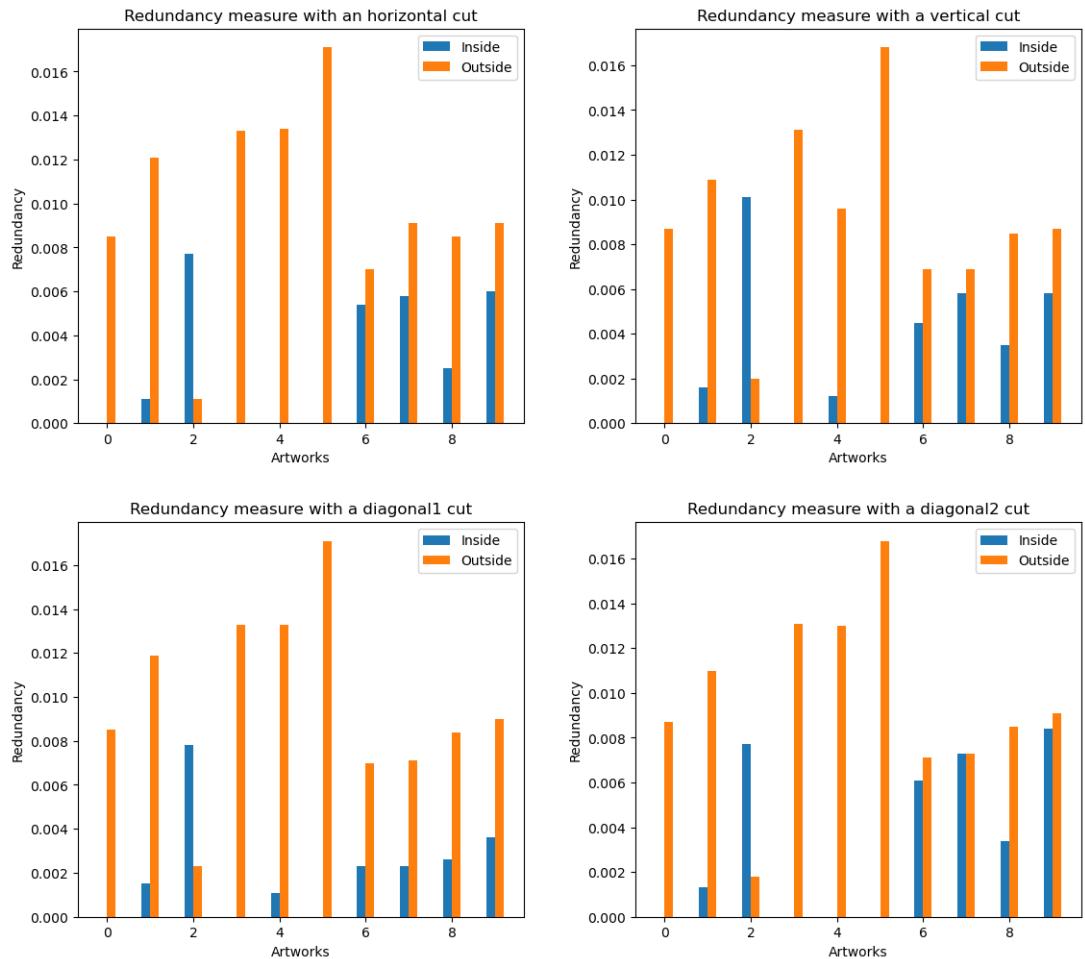


Figure 5.20: Information measures were collected for 10 distinct user triplets, yielding a total of 4 redundancy values for each triplet.

CHAPTER 5. RESULTS

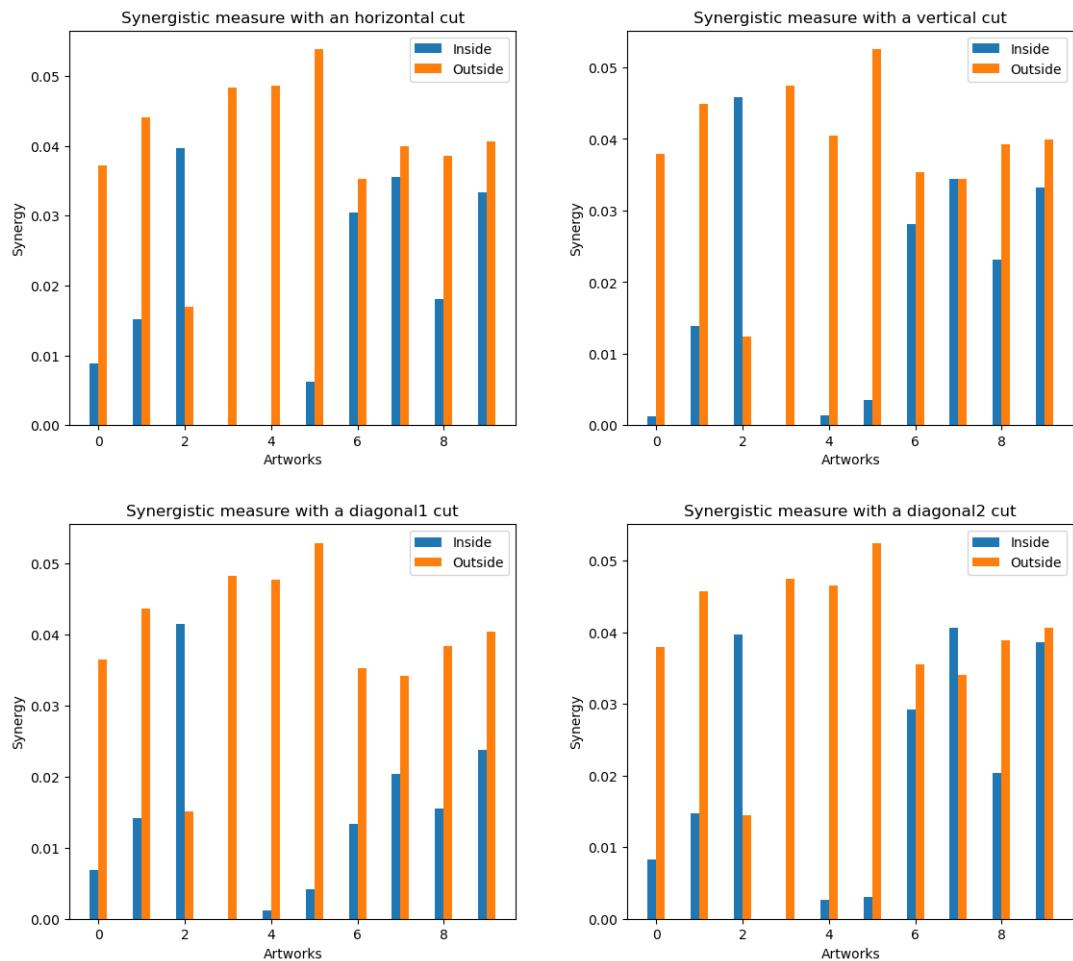


Figure 5.21: Information measures were collected for 10 distinct user triplets, yielding a total of 4 synergistic values for each triplet.

CHAPTER 5. RESULTS

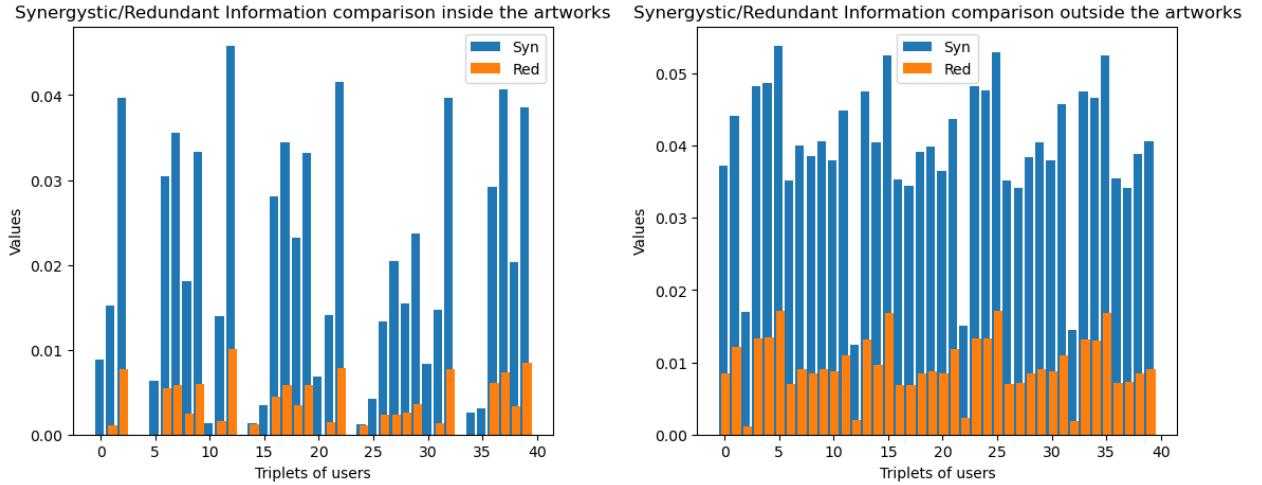


Figure 5.22: Global view of synergistic and redundancy values for the 10 different triplets of users. On the x-axis we have the triplets.

In fact, as expected previously, redundancy values refer, in some way, to the fact that users are collaborating, because they are performing actions related to the same symbol (1,2,3) of entropy. For this reason, we were expecting to see higher values of redundancy inside the artworks, imagining them active on the same zone of the selected cut in the same time interval, or probably inactive at the same time.

Intuitively one can understand that, in the way the problem is defined, synergism means that they are working on different symbols, thus different actions, in parallel. This is the reason why we expected to see synergism higher outside the artwork than inside, and in fact that is what we observe in most cases.

That said, we can avoid dividing the results by spatial cut, and put everything together. The graph (5.22) is useful to get a complete view of the results just described.

To make the analysis free from bias related to having chosen multiple users belonging to the same artwork, we performed the same measures (5.23) on 16 different triplets, each one involved in a different artwork. Results are not different from before but confirm the same trend, evidently the bias is not the cause of the values we observe.

Fig. (5.23) reflects the PED results related to the first of our interaction definitions, consequently we proceed with the results obtained for the same user triplets but with different entropy symbol definitions (with reference to subsection 4.4.3).

From now on, the definitions will include the feature color of the placed pixel; let's start with "Equal or different colors & number of times they did anything in intervals of 30 minutes", whose results are detailed in fig. (5.24).

With reference to the graph, what is immediately apparent is a sharp reversal in information values.

This time, in fact, on 16 out of 16 cases the redundancy values are higher inside the artworks than outside, while for synergy this happens 14 times out of 16, with

CHAPTER 5. RESULTS

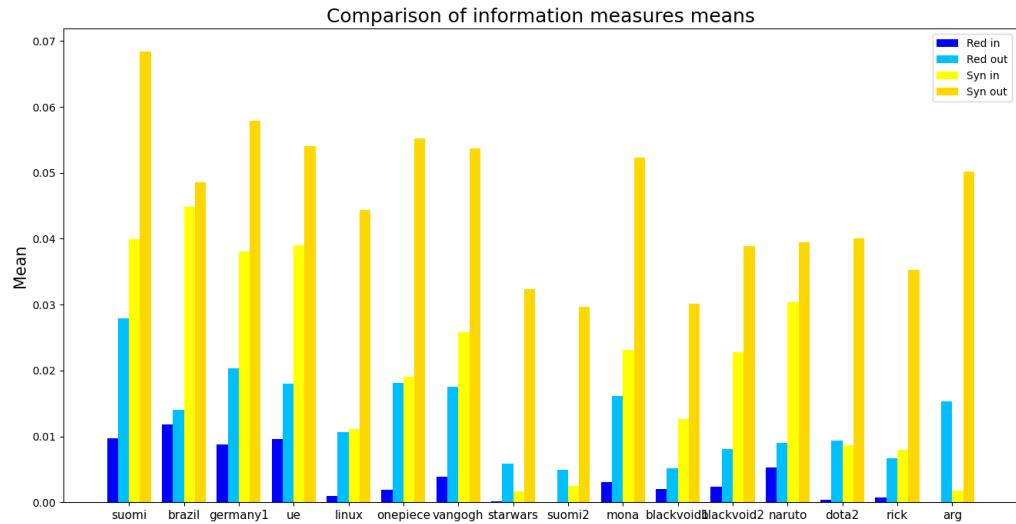


Figure 5.23: PED measures for 16 different artworks inside and outside them.

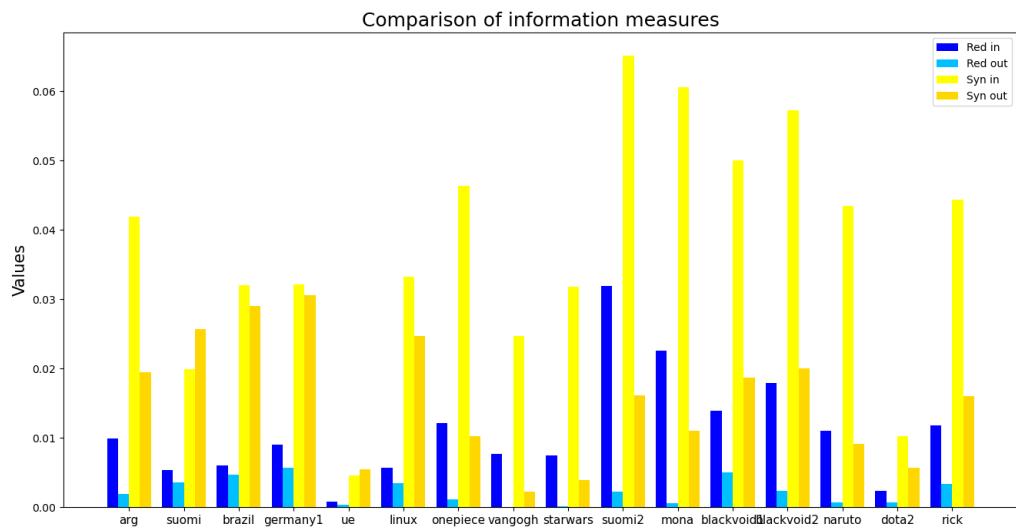


Figure 5.24: "Equal or different colors & number of times they did anything in intervals of 30 minutes" PED results

CHAPTER 5. RESULTS

very slight differences in values in the 2 cases where measures outside the artworks are higher.

It is also curious to observe that although we are looking at an output that is quite different from the previous one, the range of values of the two types of information has remained the same.

It would thus seem that the color feature, as noted already in the previous sections of the results, has considerable importance in identifying the meaning of a specific placement and thus also of a hypothetical situation of collaboration among small groups of users.

The results obtained with this definition are what we expected to find, yet it is our duty to question and fully understand the definition we are using in relation to the purpose for which we put it in place.

We are looking for patterns of collaboration among users, consequently with a preliminary idea we can think about the fact that they behave in the same way in a well-delineated space-time, but placing all pixels of the same color (like the previous definition suggests) may not necessarily have to do with "collaboration," or coordinated behavior, as we understand them.

What we need is a more stringent criterion that reflects the local situation around the individual placement but involves all features of interest to us, i.e., space, time, and color.

It is in this scenario that the majority rule, introduced before for other purposes, takes over. For this purpose we decided to try to use the majority rule with two different definitions: "Majority rule & number of times they did anything in intervals of 120 minutes", and "Majority rule & number of times they did anything in intervals of 60 minutes".

Results are visible in Fig. (5.25).

After deciding to apply the majority rule also for the PED we started our analysis first with the definition that provides 120 minutes as the time interval in which users have no activity.

This is because, having studied the activity of the "redditors" outside and inside the artwork, we knew that it was much more likely to find a high value count for the symbol '3' outside of these, both because outside them the analysis extended for more time and also because outside users had no reason to act assiduously as they would have done inside.

Therefore imposing a longer time interval (120m) where users do nothing means trying to find comparable counts inside and outside the artworks. After that we tried to reduce the timedelta to 60m to see how the output would change.

Using the 120m definition we can observe that we get the expected results, similar to the previous case shown in fig. (5.24), in which the redundancy is higher inside the artwork in 14 cases out of 16, with slight differences in the 2 cases in which this does not happen. Even the synergy is greater inside in most cases, and as usual manifests itself on values much higher than those of redundancy.

But what happens when we reduce by half the time interval of the symbol '3'? Here again the situation "worsens", in fact on 9 cases out of 16 the redundancy is greater

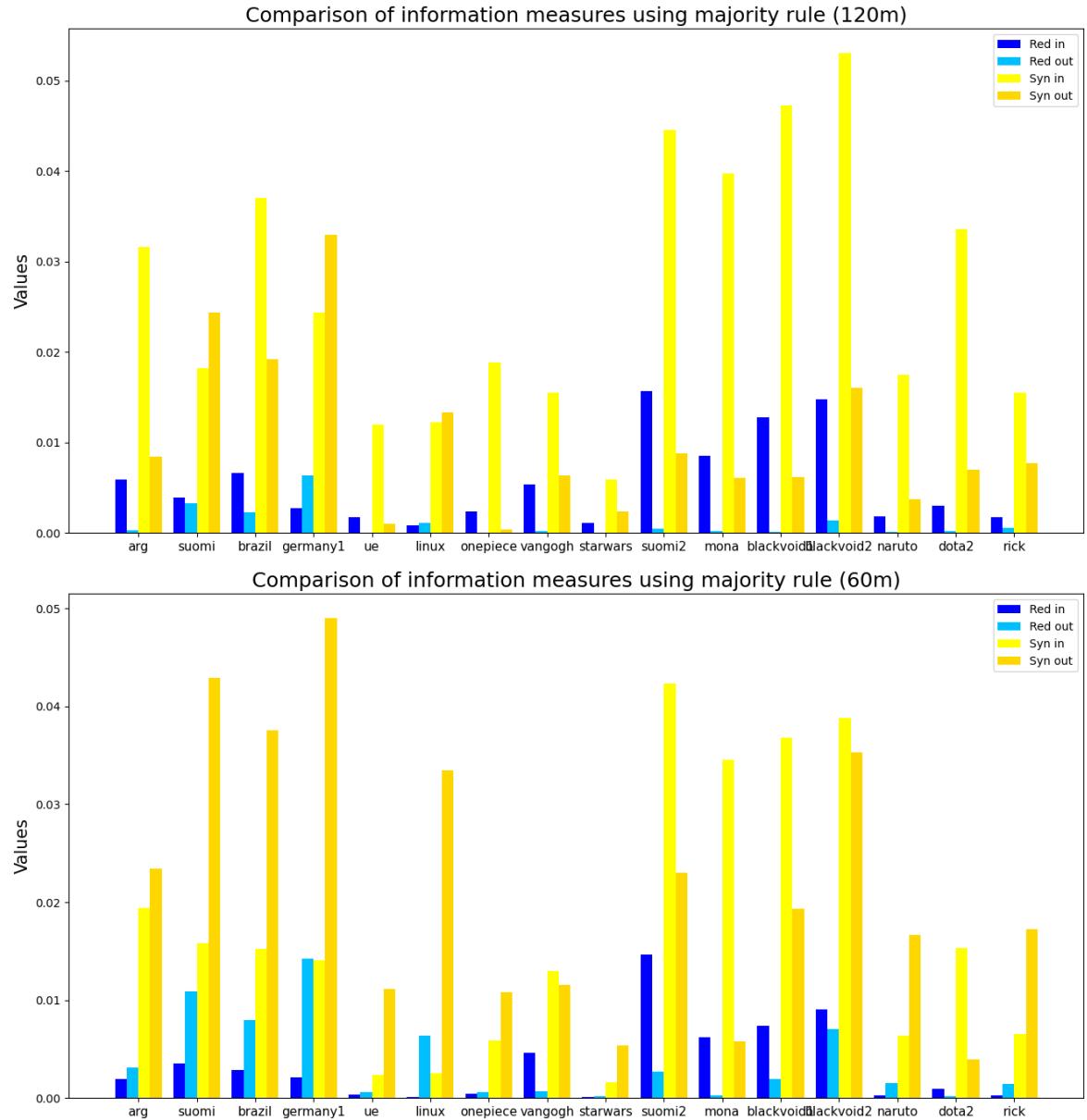


Figure 5.25: Majority rule definition's PED results. Two different timedeltas (120 and 60 minutes) selected for symbol '3'.

outside the artwork than inside, and the synergy behaves in a manner coordinated to redundancy, showing itself to be superior to that inside in the same cases in which this happens for redundancy.

So, we are now faced with two almost opposite situations, what does this mean?

First of all, we discover that one of the variables that most influence our measurements is the symbol '3', or rather what determines the output of the symbol '3', i.e., the duration of the time interval in which we do not observe activity.

Secondly, we can say that we have identified the correct timedelta for our definition, that is 120 m, because it allows to properly compare the measures we obtain in very different environments (inside and outside), allowing us to show the differences in the coordinated behavior (or not) of user triplets.

At this point, we can conclude that, trying several definitions, we found out which factors most influence the entropy measurements we perform and which information they return.

The PED is therefore a very powerful tool in the identification of behavioral patterns even in social contexts.

Based on our findings and what we've learned from previous analysis, we can conclude that the use of majority rule as a tool of search for collaboration between users through the PED is the most suitable for our purposes, showing also how the activity of users is more coordinated within the temporal space bounded by the artwork.

The PED has also been useful as a comparison tool between similar definitions, allowing us to refine the definitions of our symbols by detecting the correct timedelta for the "no-activity" symbol.

5.4.2 Toy models

In this section we describe the results obtained by performing information measures on the synthetic timeseries generated by following the procedure depicted in section 4.3.2.

First of all, we have to point out something very important about these toy models that is, having used the poissonian gillespie algorithm has led to two consequences: the first is that the duration of these generated timeseries is shorter than the original timeseries, on average, by about 13 hours; the second is that for the same reason the number of users participating in the canvas is smaller than the original timeseries by about 100000. This happens because the placements on the canvas are not Poissonian processes, but for our purposes this does not matter much, because of the easiness of the models involved.

We can move on with Fig. (5.26), who is important to have in mind the models we implemented, remembering that the parameters under study are the timedelta of interactions and the lenght of the list of the previous active users in the selected timedelta. Also we tried to order these models based on ascending realism, from the less realistic, the random one, to the most realistic, that are the ones with timedelta equal to 1 minute, as suggested by previous studies on inter event times and interactions themselves.

CHAPTER 5. RESULTS

- #1 : Each user is chosen by his rate per hour of placement and put a pixel in a **random** position;

Timedelta: 10 minutes

Each user is chosen by his rate per hour of placement and has the possibility to put a random pixel or an interacting pixel, that is a pixel whose position falls in one of the 8 surrounding position of some other user's pixel (or himself) who placed in the last 10 minutes

- #2: (10m, lenght = 10);
- #3: (10m, lenght = 100);
- #4: (10m, lenght = 1);

Timedelta: 1 minute

- #5: (1m, lenght = 10);
- #6: (1m, lenght = 1);



Timedelta: 10 seconds

- #7: (10s, lenght = 1);

Figure 5.26: Toy models summary. The parameters under study are the timedelta and the lenght of the list of the previous active users in the selected timedelta.

That said, we are ready to describe our study based on information values comparison: basically, once the timeseries has been generated, we took the same users (as triplets) of the previous analysis (hence also the same artworks) and we performed PED on them (using "Spatial cuts & number of times they did anything in intervals of 30 minutes" configuration) outside the artwork.

Of course, in the case of synthetic timeseries, won't be any artwork because the color of each placement is selected randomly, so by "outside" we mean the same geographical region outside the one in which the artwork analyzed in the actual canvas lived.

In fact, the goal is to simulate the likely random behavior of triplets of users in a space where they do not intend to cooperate in the construction of something. Consequently, the entropy values we obtain with our models will be compared with those outside the artwork in the real canvas; when a model has entropy values comparable to the real ones then that is probably an indication that the selected model is realized well to reconstructing the behavior of users outside any spatial zone.

We can, therefore, begin the description of our results: we start by comparing the measurements obtained with a random model and with the 3 models with timedelta equal to 10 minutes, with reference to fig.(5.27).

Looking at the graph, we notice that for both redundancy and synergy the 3 models with timedelta 10 minutes reach values higher than the random model in most cases, indicating us, as expected, that between the most suitable class of models to replicate the dynamics of the canvas could be the one with timedelta 10 minutes.

Otherwise, it is difficult to determine the value of the most appropriate "L" parameter for our case study, in fact the values we observe do not give us a clear indication about this. We appreciate, however, having encountered higher synergistic values of redundancy as happens in the real canvas.

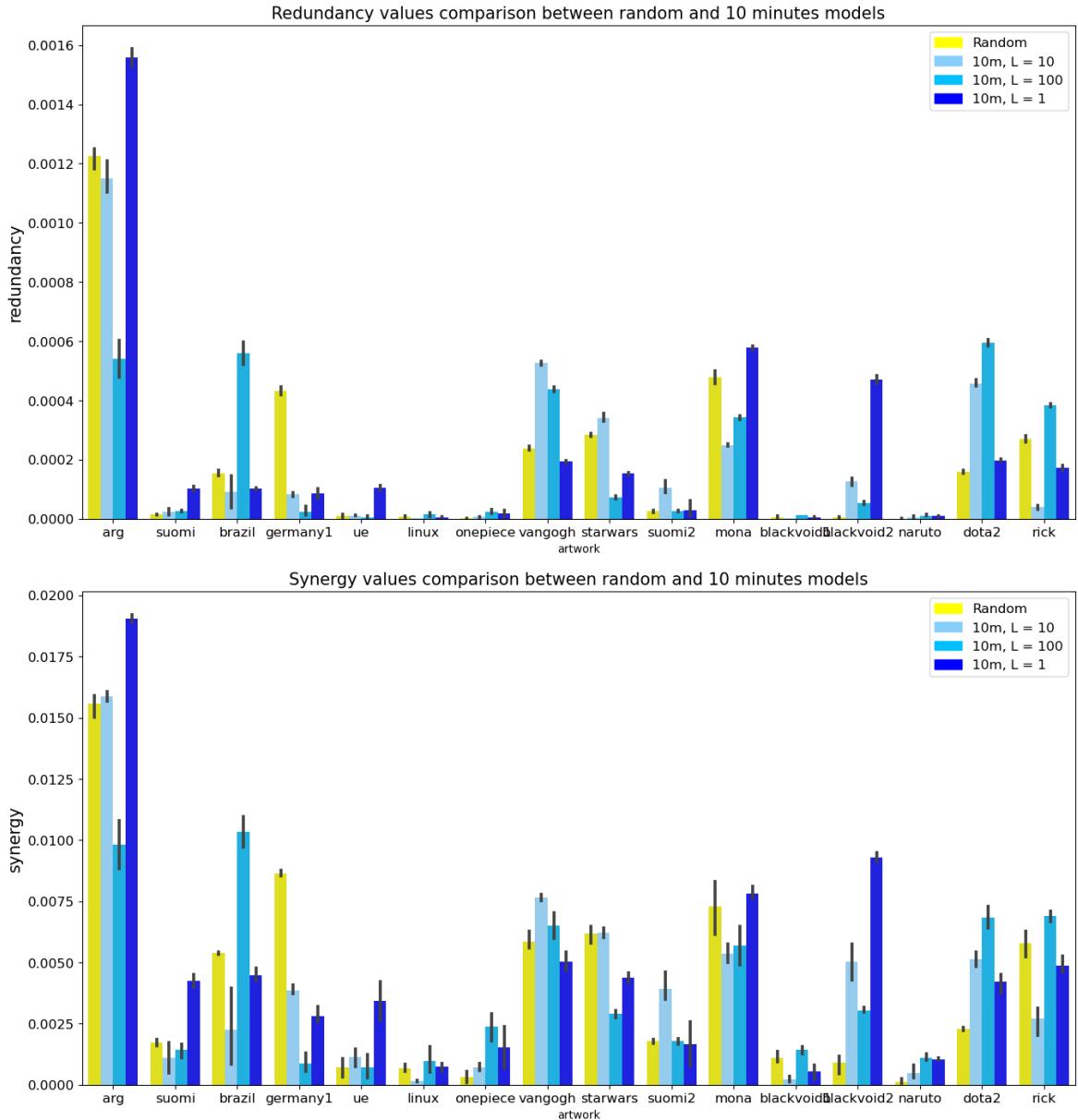


Figure 5.27: Outside measures comparison between random model and the three 10 minutes timedelta models.

CHAPTER 5. RESULTS

At this point we could go ahead showing all the other comparisons performed, i.e., 10m VS 1m models, or 10m VS 10s, but as explained in figure (5.26) we already know that a timedelta of 10 minutes is too much extended for our concept of interaction and consequently, in order of ascending realism, we put forward on it models with lower timedeltas as 1m and 10 s. For these reasons, this comparison is not graphically reported because it's redundant, even if it has been realized.

A much more interesting comparison it's definitely between 1m and 10s timedeltas models, fig.(5.28).

In this case, it seems to be difficult to find a model that perform better than the other looking at the information values encountered.

About redundancy we observe cases where the 10 seconds model takes much higher values than the other, but also cases where we have almost zero values in both models, and slightly less cases where the 1 minute model prevails. About synergism instead, the values rise from both sides, in fact the measures we observe are almost equal in many cases, and, in cases where the 1 minute model prevails, the differences are very thin.

Based on what was observed for the two different information measures, and on the fact that the 1 minute models were 3 different while the 10 second model was only one, we believe that the model that can try more to replicate the dynamics of the real canvas can be the latter.

It should also be said that through a comparison with the real data outside the artworks, whatever the class of models selected, we still get quantitative indications that apply to all the models we made. In fact, we just have to compare the model with 10 seconds timedelta with the real data.

Immediately, looking at fig.(5.29) who shows our last comparison, we realize that the values of our hypothetical best model are far from the real data in all cases. This applies to both redundancy and synergy.

In fact, the scale of values on which the real data is found has always been much greater than that of our models, since when we started the analysis with the random model.

How can we explain this findings? Surely, the approach we used is correct, what can be less so is the complexity of the models we designed. It's obvious that there is some aspect of the canvas dynamics that we don't grasp, and these measurements suggest that referring only to the spatial position of pixels in a space is not enough to identify user placement probabilities.

It is also probably true that the space outside the artwork is to be considered too large to be divided into 4 zones, from this in fact arise the probability of interaction that we need for the calculation of the PED.

Another hypothesis is that the way we built the algorithm and then the way users place pixel is also too simple to replicate the coordinated behavior of groups of three of these.

CHAPTER 5. RESULTS

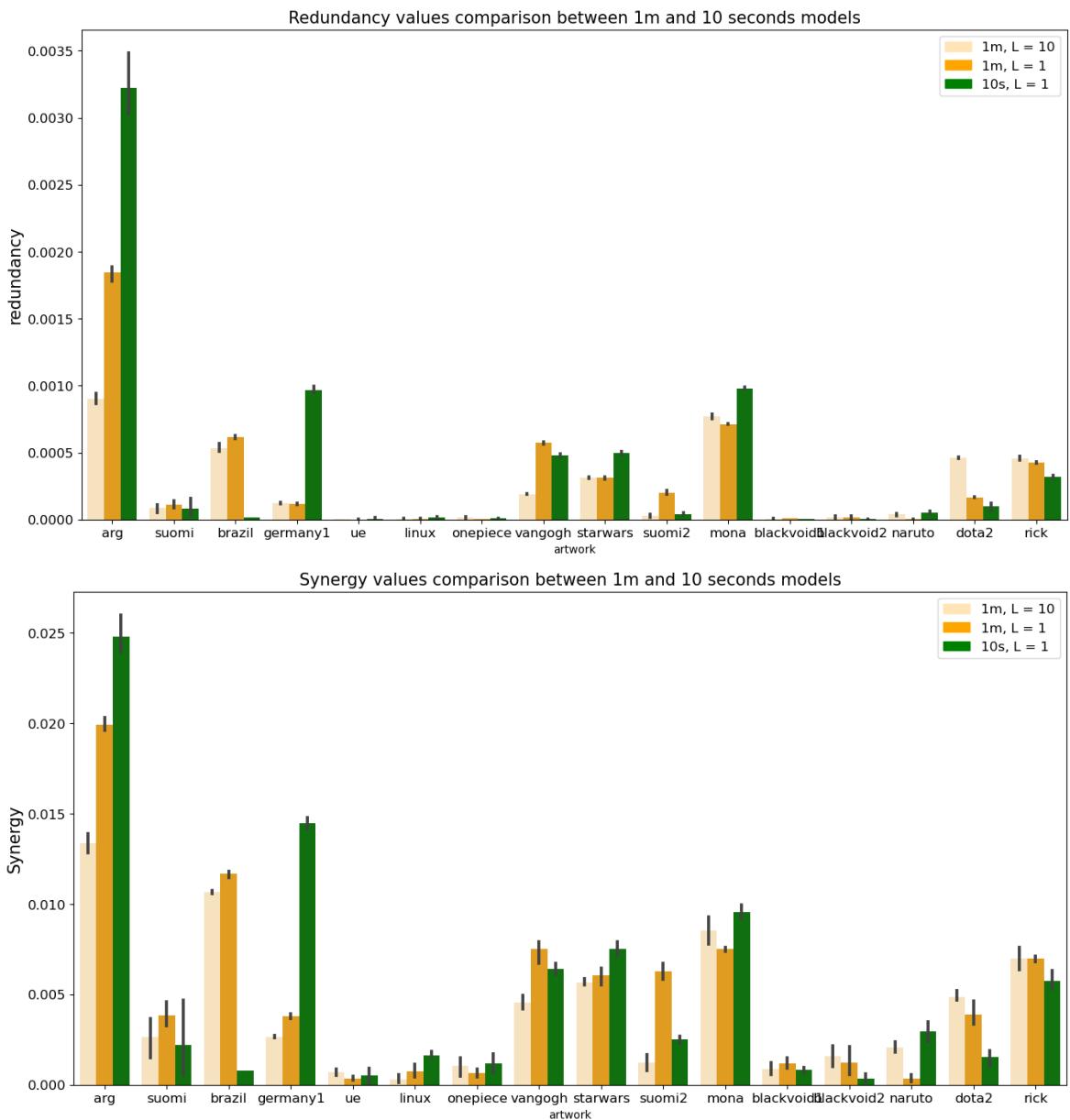


Figure 5.28: Outside measures comparison between 1 minute and 10 seconds timedelta models.

CHAPTER 5. RESULTS

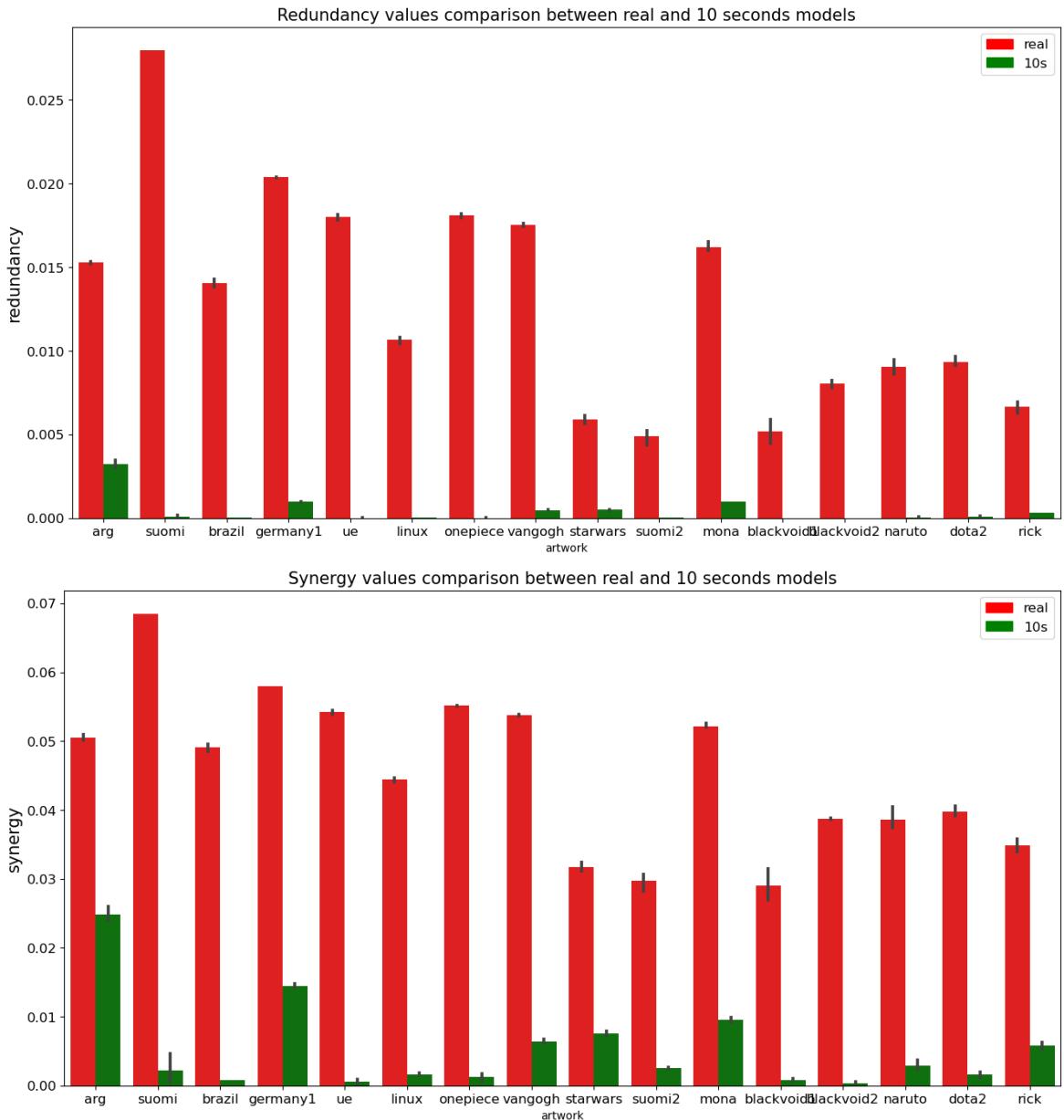


Figure 5.29: Outside measures comparison between 10 seconds model and real data.

CHAPTER 5. RESULTS

Finally, we had, however, assumed that these were not Poissonian processes, and the use of a Poissonian approach could have had a minor impact on our results.

That said, future works on this issue are definitely to be carried out, such as the use of different definitions for the PED or various timeseries generation algorithms.

Chapter 6

Discussions

In this project, our goal was to explore the characteristics of a real online social system and its resulting collaborative behavior. Our selection of the dataset was driven by this objective, and we found that the Reddit Place experiment seemed appropriate for this type of analysis, also having, in terms of social structure, an organization that easily allowed us to interpret it through high-order interactions, specifically using information theory tools, which was also part of the purpose of this research.

We started with more general analyses of the dataset and then gradually zoomed in on the problem. However, already at the global level, some interesting aspects emerged about the data, especially about users activity, such as the spatial heterogeneity of pixel placement that reveals specific trends of activity distribution across the canvas (see activity heatmap in Fig.5.1).

We also discovered that the majority of user just placed no more than 5 pixel, probably because people were curious about this social experiment but didn't have a lot of time to spend on the canvas since the presence of a 5 minutes activity-bond. Committed users, i.e. the ones who placed more than a 100 different pixel, represents, instead, a small but consolidated part of the total users in the dataset.

Talking about that, we moved on by studying inter-event time distribution in terms of subsequent pixel placed by users: was very interesting to see, also because was an great insight for the aim of this thesis, how activity of triplets of users was most intense inside selected artworks than outside them (fig. 5.6), probably indicating the presence of a coordinated behavior between them.

Was, also, useful to find out how activity of users has changed through the days of the experiment (5.8), showing up how initially the activity level was very low and gradually went up, with specific intervals when there was an activity spike. This helped us to understand at what times the users were most present.

After this preliminary analysis, we identified particular relationships between artworks, considering that each one represents a specific community who discuss by the use of a subreddit. Relations of peace and conflict came up, these allowed us to discriminate the behavior of users belonging to different communities and to see, where there were peaceful situations, if the boundaries between neighboring artworks were respected, and where there were conflicts, if users were particularly involved in wanting to eliminate the artwork of opponents.

CHAPTER 6. DISCUSSIONS

In this scenario, we devised a majority rule, designed as a tool for mapping the local meaning in terms of cooperation of a single placement but also taking into account the pixels around it. This would not have been possible without the integration of the color feature. Our findings reveal that majority rule was a strong measure to capture the specific dynamics involved between users belonging to 2 or more communities.

We, then, tried to build some definitions of interactions between users, without asserting the meaning of these, that is, without saying whether it was a collaboration or not, but with the sole meaning of asserting the fact that these users had come into contact. What needed to be understood, however, at least initially, was on what occasions these users had come into contact in a fortuitous way, due primarily to chance or acting in a limited space. For these reasons, we generated a null model by a random reshuffle of pixel positions, so we could have labeled the interactions found by this model as "random" by eliminating them from the total count of those identified in the real canvas.

We have also shown how, as the lifetime in which an interaction can exist decreases, the tendency of users goes from wanting to interact to place random position pixel, placing most of their pixels in areas that have nothing to do with a possible collaboration or some type of coordinated behavior. This also allowed us to establish that the use of a timedelta too extended, for the concept of interaction we thought, was not suitable.

This interaction analysis has been useful to answer the question about activity-interaction correlation, that is, are the most active users the ones who interact the most? The answer seemed to be negative (fig. 5.19), the only case in which we found a low-positive correlation is the case of a 30 seconds timedelta, probably we identified the most likely time interval of interaction.

Once we acquired a lot of information on our dataset, we tried to replicate the dynamics of the real canvas by generating fictitious timeseries with Gillespie's algorithm. In this case, for our concept of interaction, we tried different timedelta as the time in which the latter lives. Subsequently, to understand how our toy models behaved, we compared the information measurements of Partial entropy decomposition to specific artwork and we compared those results with those of the real canvas. The results, however, show values too dissimilar to be able to replicate what happens in the social experiment, it is likely that the models we developed are too simple to reproduce what happened on Place.

Finally, after introducing the PED, we applied it to triplets of users in specific artwork, thus generating redundancy and synergy measures aimed at representing their possible coordinated behavior. We tried different configurations for PED symbols, and applied them to triplets of users inside and outside their artworks. This allowed us to identify which entropy approach might be the most appropriate to show the presence of collaboration within the artworks and to distinguish user behavior outside of them. The majority rule, based on the fact that it takes into account both spatial proximity and the color feature, therefore seems to be the most appropriate tool to be incorporated within the PED to identify collaboration between users.

6.1 Future works

The mechanisms underlying the complex dialogue between the structure of a system and its emergent behavior such as cooperation are complex and not so obvious.

We've only just started looking into it, and there's a whole lot more we can discover about this interesting topic. We can make changes and try different approaches to learn more. Here, I'd like to mention a few things we're excited to explore further.

Starting from interactions, since we were able to build several definitions of them, one, based also on this work, could try to make new definitions of interactions unifying some of the features available from the dataset and verify their goodness.

We also think that most of the work can be done about the understanding of the dynamics of the canvas, trying to repeat the procedure we applied using the other 3 definitions of interactions we made. Another approach could be the one that uses a non-poissonian gillespie algorithm to generate the synthetic timeseries. Other ways could be followed, like using an agent based model, or any other tool trying to get closer to the real behavior of users.

Finally, although PED has proved to be a great tool for this case study, one could attempt to map user behavior in specific areas by selecting more than three users at time, even if this can lead to much greater analytical and computational efforts.

As explained in [36], it can be quite challenging to deduce higher-level structures from time series data, which capture the dynamic actions of the nodes instead of directly measured connections and relationships. Methods for reconstructing relationships solely based on temporal correlations face the limitation of not being able to completely differentiate between direct and indirect causation. In other words, they cannot distinguish whether there's a direct edge or hyperedge connecting nodes or if there's a more extended, indirect pathway connecting them. Unfortunately, these methods are unable to identify non-causal correlations.

Considering these limitations, it would be quite intriguing to revisit this dataset in the future, especially when new methods for high-order structures become available and gain substantial support.

List of Figures

| | | |
|------|---|----|
| 3.1 | Reddit place description | 12 |
| 3.2 | Reddit place artworks atlas | 14 |
| 3.3 | Some of the artworks of the final canvas | 15 |
| 4.1 | An example: brazil artwork at the beginning and at the end. | 17 |
| 4.2 | The moment when an artwork begins and ends, in this case Linux artwork | 18 |
| 4.3 | Visual representation of our definition of interaction. | 19 |
| 4.4 | PDF of inter event times per user | 21 |
| 4.5 | Redundancy lattice for (A) two variables, (B) three variables. | 25 |
| 4.6 | Full PED in the case of 2 variables | 27 |
| 5.1 | Activity heatmap of the entire experiment aggregated on a grid of 5-pixel per side for visualization. | 32 |
| 5.2 | Number of placements per user distribution | 33 |
| 5.3 | (a) Activity over time: clicks per hour (blue) and active users per hour(red). (b) Canvas coverage over time. | 34 |
| 5.4 | Inter placement times of the 10 most active users on the canvas. | 35 |
| 5.5 | On the left: PDF of inter event times of every user on the canvas. On the right: PDF of inter event times of only neighboring pixel, i.e. interacting pixel. | 35 |
| 5.6 | Boxplot distributions of inter-event times inside and outside some artworks considering the activity of one triplet of users for each one. | 36 |
| 5.7 | Activity trend in the whole experiment divided in 5min intervals. | 37 |
| 5.8 | Activity trend day by day divided in 1min intervals. | 38 |
| 5.9 | Colors distribution. | 39 |
| 5.10 | Peaceful relationships scheme between specific artworks | 40 |
| 5.11 | The conflict under our analysis: germany forced france to move from their original position. | 41 |
| 5.12 | (a) | 42 |
| 5.13 | (b) | 42 |
| 5.14 | (c) | 42 |
| 5.15 | Users behaviour analysis on different types of artworks. Suomi and Argentina were neighbors on the canvas in a peaceful relationship and this is verified by their users behaviour (a), (b). France and Germany were in a conflict, as we can see by data in (c). | 42 |

LIST OF FIGURES

| | |
|--|----|
| 5.16 Majority rule examples: on the left an example of collab pixel, on the right an example of a conflict one. | 43 |
| 5.17 PDF of rates per hour of every user | 44 |
| 5.18 Probability distribution functions of pixel per user who generated at least an interaction | 46 |
| 5.19 Activity-interactions correlation for different timedeltas | 47 |
| 5.20 Information measures were collected for 10 distinct user triplets, yielding a total of 4 redundancy values for each triplet. | 49 |
| 5.21 Information measures were collected for 10 distinct user triplets, yielding a total of 4 synergistic values for each triplet. | 50 |
| 5.22 Global view of synergistic and redundancy values for the 10 different triplets of users. On the x-axis we have the triplets. | 51 |
| 5.23 PED measures for 16 different artworks inside and outside them. | 52 |
| 5.24 "Equal or different colors & number of times they did anything in intervals of 30 minutes" PED results | 52 |
| 5.25 Majority rule definition's PED results. Two different timedeltas (120 and 60 minutes) selected for symbol '3' | 54 |
| 5.26 Toy models summary. The parameters under study are the timedelta and the lenght of the list of the previous active users in the selected timedelta. | 56 |
| 5.27 Outside measures comparison between random model and the three 10 minutes timedelta models. | 57 |
| 5.28 Outside measures comparison between 1 minute and 10 seconds timedelta models. | 59 |
| 5.29 Outside measures comparison between 10 seconds model and real data. | 60 |

List of Tables

| | | |
|-----|--|----|
| 5.1 | Users analysis in different types of artworks: argentina and suomi reached an agreement of peace, instead of France who was involved in a conflict with Germany. At the end also Germany lost his position, transforming itself into EU artwork. | 43 |
| 5.2 | Interactions analysis for different timedeltas on the original canvas. Random interactions refer to the number of interactions we found using the respective null model. | 45 |

Bibliography

- [1] L. Alessandretti, U. Aslak, and S. Lehmann, “The scales of human mobility,” *Nature*, vol. 587, pp. 402–407, nov 2020.
- [2] C. Liu and Z.-K. Zhang, “Information spreading on dynamic social networks,” *Communications in Nonlinear Science and Numerical Simulation*, vol. 19, no. 4, pp. 896–904, 2014.
- [3] A. Bazghandi, “Techniques, advantages and problems of agent based modeling for traffic simulation,” 2012.
- [4] A. C. Cullen, T. Alpcan, and A. C. Kalloniatis, “Adversarial decisions on complex dynamical systems using game theory,” *Physica A: Statistical Mechanics and its Applications*, vol. 594, p. 126998, 2022.
- [5] T. F. Varley, “Information theory for complex systems scientists,” 2023.
- [6] T. F. Varley, M. Pope, null, null, and O. Sporns, “Partial entropy decomposition reveals higher-order information structures in human brain activity,” *Proceedings of the National Academy of Sciences*, vol. 120, no. 30, p. e2300888120, 2023.
- [7] J. Andreas, G. Beguš, M. M. Bronstein, R. Diamant, D. Delaney, S. Gero, S. Goldwasser, D. F. Gruber, S. de Haas, P. Malkin, N. Pavlov, R. Payne, G. Petri, D. Rus, P. Sharma, D. Tchernov, P. Tønnesen, A. Torralba, D. Vogt, and R. J. Wood, “Toward understanding the communication in sperm whales,” *iScience*, vol. 25, no. 6, p. 104393, 2022.
- [8] B. Zhang and D. L. DeAngelis, “An overview of agent-based models in plant biology and ecology,” *Annals of Botany*, vol. 126, no. 4, pp. 539–557, 2020.
- [9] F. Lorig, E. Johansson, and P. Davidsson, “Agent-based social simulation of the covid-19 pandemic: A systematic review,” *JASSS: Journal of Artificial Societies and Social Simulation*, vol. 24, no. 3, 2021.
- [10] R. Cohen, S. Havlin, and D. ben Avraham, “Efficient immunization strategies for computer networks and populations,” *Physical Review Letters*, vol. 91, dec 2003.
- [11] C. Kasper, M. Vierbuchen, U. Ernst, S. Fischer, R. Radersma, A. Raulo, F. Cunha-Saraiva, M. Wu, K. B. Mobley, and B. Taborsky, “Genetics and developmental biology of cooperation,” *Molecular ecology*, vol. 26, no. 17, pp. 4364–4377, 2017.

BIBLIOGRAPHY

- [12] D. M. Kreps, P. Milgrom, J. Roberts, and R. Wilson, “Rational cooperation in the finitely repeated prisoners’ dilemma,” *Journal of Economic theory*, vol. 27, no. 2, pp. 245–252, 1982.
- [13] D. Zhou, J. Huang, and B. Schölkopf, “Learning with hypergraphs: Clustering, classification, and embedding,” *Advances in neural information processing systems*, vol. 19, 2006.
- [14] W. Powell, D. White, K. Koput, and J. Owen-Smith, “Network dynamics and field evolution: The growth of interorganizational collaboration in the life sciences,” *American Journal of Sociology*, vol. 110, no. 4, pp. 1132–1205, 2005.
- [15] B. Cooper, A. E. Lewis-Pye, A. Li, Y. Pan, and X. Yong, “Establishing social cooperation: The role of hubs and community structure,” *Network Science*, vol. 6, no. 2, pp. 251–264, 2018.
- [16] A. Li, L. Zhou, Q. Su, S. P. Cornelius, Y.-Y. Liu, L. Wang, and S. A. Levin, “Evolution of cooperation on temporal networks,” *Nature communications*, vol. 11, no. 1, p. 2259, 2020.
- [17] J. L. Juul, A. R. Benson, and J. Kleinberg, “Hypergraph patterns and collaboration structure,” *arXiv preprint arXiv:2210.02163*, 2022.
- [18] L. G. Mojica, “Modeling trolling in social media conversations,” 2016.
- [19] S. Kumar, W. L. Hamilton, J. Leskovec, and D. Jurafsky, “Community interaction and conflict on the web,” in *Proceedings of the 2018 World Wide Web Conference on World Wide Web - WWW ’18*, ACM Press, 2018.
- [20] H. Fang, H. Cheng, and M. Ostendorf, “Learning latent local conversation modes for predicting community endorsement in online discussions,” 2016.
- [21] B. D. Horne and S. Adali, “The impact of crowds on news engagement: A reddit case study,” 2017.
- [22] C. Tan and L. Lee, “All who wander,” in *Proceedings of the 24th International Conference on World Wide Web*, International World Wide Web Conferences Steering Committee, may 2015.
- [23] A. N. Medvedev, R. Lambiotte, and J.-C. Delvenne, “The anatomy of reddit: An overview of academic research,” *Dynamics On and Of Complex Networks III: Machine Learning and Statistical Physics Approaches 10*, pp. 183–204, 2019.
- [24] Muller and J. Winters, “Compression in cultural evolution: Homogeneity and structure in the emergence and evolution of a large-scale online collaborative art project,” *Twelfth International AAAI Conference on Web and Social Media*, 2018.
- [25] P. Bromiley, N. Thacker, and E. Bouhova-Thacker, “Shannon entropy, renyi entropy, and information,” *Statistics and Inf. Series (2004-004)*, vol. 9, pp. 2–8, 2004.

BIBLIOGRAPHY

- [26] B. Armstrong, “Coordination in a peer production platform: A study of reddit’s /r/place experiment,” *UWSpace*, 2018.
- [27] J. Rappaz, M. Catasta, R. West, and K. Aberer, “Latent structure in collaboration: The case of reddit r/place,” *Proceedings of the International AAAI Conference on Web and Social Media*, vol. 12, Jun. 2018.
- [28] S. C. R. Elías Gabriel Gil, “La inteligencia colectiva en la producción cultural. análisis del experimento social “place” en la plataforma web de reddit,” 2020.
- [29] P. Vachher, Z. Levonian, H.-F. Cheng, and S. Yarosh, “Understanding community-level conflicts through reddit r/place,” in *Conference Companion Publication of the 2020 on Computer Supported Cooperative Work and Social Computing*, CSCW ’20 Companion, (New York, NY, USA), p. 401–405, Association for Computing Machinery, 2020.
- [30] K. T. Litherland and A. I. Mørch, “Instruction vs. emergence on r/place: Understanding the growth and control of evolving artifacts in mass collaboration,” *Computers in Human Behavior*, vol. 122, p. 106845, 2021.
- [31] G. P. Liber Dorizzi, “Emergence of cooperation from group interaction patterns,” 2022.
- [32] A. M. Adams, J. Fernandez, and O. Witkowski, “Two ways of understanding social dynamics: Analyzing the predictability of emergence of objects in reddit r/place dependent on locality in space and time,” 2022.
- [33] D. T. Gillespie, “A general method for numerically simulating the stochastic time evolution of coupled chemical reactions,” *Journal of Computational Physics*, vol. 22, no. 4, pp. 403–434, 1976.
- [34] P. L. Williams and R. D. Beer, “Nonnegative decomposition of multivariate information,” 2010.
- [35] R. A. A. Ince, “The partial entropy decomposition: Decomposing multivariate entropy and mutual information via pointwise common surprisal,” 2017.
- [36] F. Battiston, E. Amico, A. Barrat, G. Bianconi, G. F. de Arruda, B. Franceschiello, I. Iacopini, S. Kéfi, V. Latora, Y. Moreno, M. M. Murray, T. P. Peixoto, F. Vaccarino, and G. Petri, “The physics of higher-order interactions in complex systems,” *Nature Physics*, vol. 17, pp. 1093–1098, oct 2021.

Ringraziamenti

Un doveroso ringraziamento va ai miei relatori, Giovanni e Michele, i quali ringrazio molto per avermi permesso di vivere l'esperienza della tesi all'estero. L'aver condiviso con voi scambi di idee e di pensiero è stato estremamente formativo, porterò con me i momenti di confronto che abbiamo avuto, vi ringrazio per la vostra disponibilità nonostante i tanti impegni.

Un altro ringraziamento va a Davide del gruppo di ricerca NPLab, il quale è stato preziosissimo nel chiarire i miei dubbi e nel proporre nuovi approcci per affrontare alcuni temi di questa tesi. Grazie anche a Liber, il quale si è dimostrato fin da subito bendisposto a condividere e a spiegare alcuni dei risultati ottenuti con la sua tesi, permettendomi di avere delle solide basi da cui partire per la mia.

Gracias a los chicos de la UPC, amables y buenisima gente, no garantizo nada pero intentaré mejorar mi nivel de catalan.

Vorrei tanto ringraziare Andrea A. per avermi aperto le porte del mondo che con tanta fatica si era creato in quel di Barcellona, ti sono riconoscente, mi dispiace che la mia esuberanza ti abbia quasi fatto perdere il lavoro.

Gracias Blanqueria 13, y especialmente gracias Luz, me trataste como a un hermano y eso no se daba por sentado. Espero haber correspondido al menos parte de tu cariño con una buena moka italiana.

Rigrazio gli amici di sù, per avermi accompagnato dal giorno zero in questo percorso, per aver condiviso con me lacrime e sorrisi, aule studio e discoteche, aperitivi in presenza e da remoto, in generale per esserci stati. Ve ne sono grato.

Rigrazio gli amici di giù, perchè anche quando ognuno ha preso la sua strada i rapporti non sono cambiati. Le risate per le solite demenzialità che solo noi capiamo sono state quello che mi ha sollevato quando ne ho avuto bisogno. Non crescite mai.

Il ringraziamento più grande non può che essere alla mia famiglia. Grazie mille a tutti i miei nonni, avete contribuito anche a voi a questo mio successo personale, dai bei ricordi che mi ha lasciato chi oggi non c'è più ai "non ti dico niente" di chi per fortuna è ancora qua a fare il tifo per me, siete stati fondamentali.

Grazie a tutti i miei zii per l'attenzione e il costante interesse ai miei piccoli traghetti, ma anche per tutto il vostro affetto e per tutte le vostre "m'berete".

Grazie Carlotta per il tuo incessante volermi bene nonostante la distanza. Per quanto sia grande il tuo astio verso la fisica, ora potrai vantarti di avere un fratello che ti ha mostrato come non sia poi così impossibile riuscire a capirla, o almeno, a far finta di capirla.

BIBLIOGRAPHY

Un enorme grazie a Davide, c'è tanto di te in questo mio piccolo traguardo. Il bene che mi vuoi è inquantificabile, sono fiero di averti avuto al mio fianco nelle situazioni più disparate. La tua presenza non è mai mancata anche quando credevo di non averne bisogno, sei per me la definizione di fratello maggiore.

Mamma e Papà, un semplice grazie non basterebbe per potervi dire quanto vi sono riconoscente, per me siete stati imprenditori, ultras, life coaches e tanto altro. Il vostro sostegno e amore incondizionato è qualcosa che nemmeno le leggi della fisica possono spiegare, almeno per adesso ("occhiolino, occhiolino"). Con questo piccolo successo provo a restituirlvi una piccola parte della fiducia che avete riposto in me. Vi voglio un bene matto.

Infine, Luana, a ciò che sei stata per me durante questi anni non trovo una definizione che ti renda giustizia. Grazie della pazienza che hai avuto, grazie del tuo costante sostegno, grazie di avermi rimesso in piedi quando ho sentito di non essere all'altezza. Nonostante i molti momenti passati spazialmente lontani, non ho mai provato distanza sentimentale ma solo vicinanza, felicità e conforto, ed è merito tuo. Posso solo dire di essere fortunato.

E grazie a me, perchè me lo sono meritato.