# The Streetlight Effect in Data-Driven Exploration *

Johannes Hoelzemann
University of Vienna

Gustavo Manso
UC Berkeley

Abhishek Nagaraj
UC Berkeley & NBER

Matteo Tranchero
UC Berkeley

April 21, 2024

## Abstract

We examine innovative contexts like scientific research or technical R&D where agents must search across many potential projects of varying and uncertain returns. Is it better to possess incomplete but accurate data on the value of some projects, or might there be cases where it is better to explore on a blank slate? While more data usually improves welfare, we present a theoretical framework to understand how it can unexpectedly decrease it. In our model of the streetlight effect, we predict that when data shines a light on attractive but not optimal projects, it can severely narrow the breadth of exploration and lower individual and group payoffs. We test our predictions in an online lab experiment and show that the availability of data on the true value of one project can lower individual payoffs by 17% and reduce the likelihood of discovering the optimal outcome by 54% compared to cases where no data is provided. Suggestive empirical evidence from genetics research illustrates our framework in a real-world setting: data on moderately promising genetic targets delays valuable discoveries by 1.6 years on average. Our paper provides the first systematic examination of the streetlight effect, outlining the conditions under which data leads agents to look under the lamppost rather than engage in socially beneficial exploration.

**JEL Classification:** C73, C92, D81, D83.
**Keywords:** Data; Innovation; Exploration and exploitation; Streetlight effect; Dynamic public-good problem; Laboratory experiment; Genetics research.

---

# 1  Introduction

A key scientific challenge in medical research is to sift through over 19,000 gene candidates to discover the genetic drivers of human disease. Puzzlingly, even after over two decades since the Human Genome Project mapped the entire genome, the genetic space appears relatively unexplored (Edwards et al., 2011; Nguyen et al., 2017). Only 10% of all genes have been targeted by approved drugs despite significant investment in this area and the awareness that better genetic targets might exist among less studied genes (Gates et al., 2021; Stoeger et al., 2018). These patterns are mirrored in a variety of related settings such as venture capital or industrial R&D where agents should have strong incentives to explore widely, and yet, it seems like many potentially valuable options are not pursued and collective exploration is limited. Examining the drivers of such underexploration is critical because we live in an era of apparent diminishing returns to research effort (Jones, 2009; Gordon, 2016; Bloom et al., 2020). The extent to which such diminishing returns are due to constrained search in pre-existing domains as compared to true limits to the returns from innovative activity remains an open question.

In this paper, we develop a simple theoretical framework and provide empirical evidence to shed light on this phenomenon. Our starting point is the observation that search activity does not proceed on a blank slate. For instance, as is true in most domains, a scientist beginning to explore a gene-disease link has access to data on past attempts in this area. We develop a framework to understand how such data on the performance of past exploration guides the direction of future investment. We are motivated by the parable of the *streetlight effect*, where knowledge of past discoveries leads agents to narrow search for reasons of data availability rather than adopt a wider search aperture based on reasons of market size or policy importance. In our model, we show how information on past discoveries can constrain search and, ironically, lower individual and social outcomes. Such an effect starkly contrasts with the view that (accurate) data provision can only be helpful because it reduces uncertainty and makes exploration more efficient. Our paper tries to reconcile both points of view by studying how the streetlight effect might emerge in exploratory search among rational agents and outlining the conditions under which data hampers rather than spurs individual and social outcomes.

Our questions are even more relevant as big data and artificial intelligence (AI) technologies are touted to significantly increase the productivity of innovative activity (Cockburn et al., 2019). These technologies are enabled by pre-existing data since they learn from past discoveries to provide predictions on the viability of new investments (Kim, 2023; Tranchero, 2024). For example, venture capitalists and real estate investors use data and analytics on firm performance to decide where to invest (Ewens

et al., 2018; Raymond, 2024), pharmaceutical firms leverage genetic data to guide drug development (Kao, 2024), and exploration firms exploit satellite images to decide where to look for the next mineral deposit (Nagaraj, 2022). If these technologies magnify the distortionary effects of past data on future exploration, then their effects on innovation might not be entirely positive. Understanding the conditions under which data-driven technologies could lead to suboptimal outcomes for both individual performance and social welfare thus becomes even more crucial.

We begin by sketching a simple theory of how the availability of data might shape exploration choices and individual and social outcomes. In our strategic multi-armed bandit framework, agents choose among risky projects over two periods and both actions and payoffs are perfectly observable. Risky projects can be either of low, medium, or high value, but their quality can only be learned by exploring them. In each period, the decision-maker exploits the information they already have to decide whether to invest in a pre-explored project or to take a risk by exploring a new project. When they explore, they not only bear the risks of exploration but also generate new data on the value of the project for others. In this setup, we examine how providing data on the value of one opportunity can shape exploration outcomes. Our key result is that the effects of data provision depend crucially on the type of project illuminated: data that sheds light on a medium-value opportunity can *reduce* individual and group payoffs relative to not having data on the value of any project, while data on low and high-value opportunities benefit agents and society.

The intuition behind this result is that whenever the value of investing in a risk-free medium project is higher than the expected value of other risky projects, it might become individually rational for all agents to choose the medium-value opportunities highlighted by the data. This induces herding behavior, where the data acts like a magnet and kills the diversity of search efforts. By reducing the exploration of new alternatives, data on past experiments endogenously crowd out new data generation, which harms private and collective welfare in the long run. Absent data, agents might initially make low-payoff choices but are more likely to learn from a broader exploration of the choice landscape and choose the optimal project in later periods. Using a bandit framework to depict collective search is helpful in highlighting the economic force that gives rise to the streetlight effect among rational agents: lack of coordination stemming from free-riding behavior. The positive informational externality associated with a player's exploration decision gives rise to a public good problem in the form of dynamically evolving information about the agents' common state of the world. We show how, depending on the space of parameters, this can prevent coordination and give rise to the streetlight

effect.

While our theoretical framework raises an interesting hypothesis, it is still an open question whether it is consistent with realized behavior. Accordingly, we developed and implemented an online lab experiment to test our predictions. In the experiment, groups of players collectively engaged in a two-period game of strategic exploration. They were each presented with a choice of five options with unknown but varying payoffs drawn from a known distribution. In the first period, participants were instructed to sequentially choose one project whose value would only be revealed after everyone chose. In the second period, they could see the payoffs of the project chosen in the first period by all participants before making their new choice. Payoffs were non-rival and cumulative: players earned the full sum of payoffs from their choices across the two periods, irrespective of whether other players made the same choice. In this setup, participants in the baseline condition were not shown any information, while those in other conditions were shown the payoff of one (low, medium, or high value) project at the outset of the game.

The experimental results are consistent with the predictions of our framework. Results show that data provision on the medium-value project reduces individual payoffs by roughly 17% and reduces the likelihood of finding the optimal outcome by 54% compared to the condition without any initial data. In other words, given a choice, rational agents would rather explore in the dark than be told where the medium value project is. On the other hand, data describing either the low or the high-value projects raise both individual and group payoffs. In line with our theoretical framework, the mechanism is that data on medium-value projects lowers early-stage exploration and reduces endogenous data generation at the group level. We observe free-riding on data even when the payoffs are such that coordination is a subgame perfect equilibrium, suggesting that strategic concerns can prevent the emergence of coordinated exploration in parallel search. We also investigate why some players choose to explore against their own interests, ruling out that it is due to learning, imperfect understanding of the game, or risk-loving attitudes. Qualitative evidence from participants' answers suggests that deviations from the optimal behavior are driven by occasional gambles or exploration of unknown options to avoid boredom. Overall, the experimental evidence confirms that such deviations are not enough to offset the free-riding behavior that gives rise to the streetlight effect.

Finally, we empirically explore if the dynamics we highlight are consistent with patterns observed in high-stakes, real-world settings. Specifically, we return to our opening example around genetics research and map our existing framework to rich data coming from DisGeNET, a unique database

that links scientific publications with the disease and genes studied. Importantly, DisGeNET assigns a stable score to each association between a gene and a disease, with higher scores indicating greater evidence of a scientifically valuable discovery. We use this normalized score to define which genetic discovery has a low, medium, or high scientific promise for a given disease. The final dataset includes the complete information on genetic discovery for 4,369 diseases over the period 1980-2019, along with a measure of scientific value for each gene-disease association, providing us with a complete description of search activity in this crucial field.

Our conceptualization of the streetlight effect would predict that scientists exploring without the guidance of any past data would be better off than those exploring in the presence of data suggesting moderately promising gene-disease links. This is precisely what we find. Diseases for which a gene of medium value was identified early on took +7.2% more time to find the gene behind the disease as compared to diseases where no promising candidate was initially found. This difference amounts to an additional delay of about 1.6 years on average. In line with our theory, we show that the reason for this counterintuitive result is that early hits ended up reducing the diversity of follow-on research, leading to lower exploration of new gene-disease links. While these results are suggestive rather than causal, they can be visualized with clean event-study specifications that document the decline in the diversity of follow-on exploration after discovering medium-value genes. We also confirm the result with extensive robustness tests, such as accounting for total effort and alternative definitions of what constitutes a medium or high-promise genetic association. Taken together, the empirical evidence from genetics provides striking evidence consistent with our proposed theory.

Our three-part study contributes to several strands of research. First and foremost, we add to a nascent literature on the nature of data, how they are generated, and how they shape economic outcomes (Bessen et al., 2022; Farboodi and Veldkamp, 2020; Jones and Tonetti, 2020; Nagaraj and Stern, 2020; Bergeaud et al., 2024; Bergemann and Bonatti, 2019). Instead of considering data as a homogeneous commodity, we show that the nature of the data itself (in particular what it does and does not highlight) shapes agents' exploration choices. Notably, our results emerge in a context where we operationalized data as instrumental information, i.e., unbiased and directly payoff-relevant. Our results could be even starker if data were imprecise or biased (Henrich et al., 2010; Cao et al., 2023) or if agents' attention is drawn to salient payoffs (Bordalo et al., 2012, 2013, 2020). Additionally, we also propose a novel mechanism through which data might harm exploration. Our theoretical framework shows how data can cause agents to implicitly coordinate on certain – but dominated – projects and thus lower overall

exploration activity, harming group and individual outcomes by crowding out new data generation.

Second, we speak to the literature on strategic experimentation and social learning (Schlag, 1998; Bolton and Harris, 1999; Keller et al., 2005; Klein and Rady, 2011; Hörner et al., 2022). We add to recent advances exploring strategic interdependencies among players (Boyce et al., 2016; Hoelzemann and Klein, 2021, 2024), showing how informational spillovers in collective experimentation settings can give rise to free-rider problems that endogenously curtail exploration and thus aggregate data generation. Furthermore, we find that this is consistent with the patterns observed in scientific research on the genes that cause human diseases (Gates et al., 2021; Edwards et al., 2011; Haynes et al., 2018), showing how our framework can help rationalize the dynamics observed in observational data from real search tasks.

Finally, we build on the innovation search literature that has explored the drivers of risky search among innovative agents (March, 1991; Levinthal, 1997; Manso, 2011; Azoulay et al., 2011; Ederer and Manso, 2013; Henry et al., 2022). We highlight the role of the information environment in driving underexploration. We also contribute to research studying the importance of different types of data in shaping experimentation decisions in risky environments (Ewens et al., 2018; Camuffo et al., 2020; Krieger, 2021). In particular, we show how data might have counterintuitive effects in search, offering a less sanguine model for how innovation will be shaped in the age of big data and AI (Kim, 2023; Ajay et al., 2018; Cockburn et al., 2019).

The remainder of the paper proceeds as follows. In Section 2, we provide an overview of the theoretical framework, including a simplified formal model and a numerical illustration. Sections 3 and 4 describe the experimental setup and results, respectively. Section 5 presents the empirical analysis in the context of genetic research. Section 6 concludes.

## 2   Theoretical Framework

**Setup**   There are $n$ agents engaged in a search to maximize their individual payoffs and must choose between projects of initially unknown value. For simplicity, consider a scenario with 5 projects and $n \geq 5$. Such projects can have three types of payoffs: 3 have a low payoff $L$, 1 has a medium payoff $M$, and the remaining 1 has a high payoff $H$, such that $0 < L < M < H$. This distribution is known ex ante to the agents, but they do not know the type of any specific project at the outset of the game. All agents live for two periods and are risk-neutral with zero discounting. Agents cannot communicate

with each other. This setup represents many settings where agents are faced with a finite number of options of unknown value, where "good" projects are hard to find but have high payoffs (Kerr et al., 2014; Manso, 2016).

**Dynamics**  The game unfolds as follows. The $n$ agents sequentially choose a project in each period according to a random order. They can observe the options that players who moved earlier selected, but do not yet learn the value associated with their choice. Once all agents have selected a project, the underlying payoffs of their chosen project are revealed to all players, and period 1 concludes. In period 2, agents repeat this process, knowing the payoffs associated with previously explored projects. Similar to period 1, agents choose sequentially according to the same random order. They can select a previously explored project of known value or an unexplored one, the value of which will be revealed at the end of the second period. Once all choices are made, period 2 ends. Payoffs are cumulative across the two periods, i.e., the sum of values associated with their choices over time, and projects are assumed to be non-rival. If multiple agents choose the same project, they all receive its payoff. This setup mimics competitive markets where organizations engage in parallel research and development: projects do not directly compete, but the generation of information about what works and what does not is valuable for all participants in the market (Krieger, 2021). In contrast to conventional payoff externalities in public good problems, the presence of the other agents impacts a given agent only via the data they produce over time (Hoelzemann and Klein, 2021, 2024).

**Equilibrium without Data**  To set the stage, we lay out the setting in the absence of any data, meaning that no information about payoffs is disclosed to agents at the outset of the game.

**Proposition 1.** *The equilibrium without data involves all projects being explored in period $1$ and the high-value project (breakthrough) being selected by all agents in period $2$. The expected payoff to each agent is: $\frac{3}{5}L + \frac{1}{5}M + \frac{6}{5}H$.*

In this simple setup, agents are initially indifferent between choosing projects since they all have the same expected payoff. However, the sequential order of choice guarantees that rational agents select a different project than the ones already chosen, in order to reveal more information that might be useful in the second period. Said otherwise, sequential choosing permits implicit coordination and ensures that all projects will be explored in the first period as $n \geq 5$. The expected payoff is the likelihood of a random draw in period $1$, and of $H$ in period $2$. To summarize, dispersed exploration in the absence of starting data ensures a breakthrough in the first period – intended as the discovery of the maximum – and high payoffs in the second period.

**Equilibrium with Data on $L$ or $H$ Projects**   We compare the setup outlined above with a setting where data about the payoff of one project is publicly revealed at the outset of the game. Depending on which project is disclosed, different dynamics unfold. We begin with the two cases where data either rule in the best alternative or rule out a poor option. Let $\pi_i$ denote a player's payoff where $i \in \{\emptyset, L, M, H\}$ indicates data provided and $P(H|i)$ be the conditional probability of discovering $H$ given data $i$.

**Proposition 2.** *If the underlying value of a project is revealed to be high, every agent selects the project whose value is revealed and achieves a payoff of $2H$. If the project revealed is low value, the data rules out one dominated option and the payoff expected by each agent is $\frac{2}{4}L + \frac{1}{4}M + \frac{5}{4}H$. Taken together, we have the following two results:*

$$[\text{Breakthrough}] \quad P(H|\emptyset) = P(H|L) = P(H|H) = 1 \tag{1}$$

$$[\text{Payoffs}] \quad \pi_\emptyset < \pi_L < \pi_H \tag{2}$$

If the $H$ project is revealed, all underlying uncertainty is resolved and each player selects that particular project in both periods. This highlights how data can guide discovery by leading directly to the best outcome (Nagaraj, 2022). If the project revealed returns a low payoff, then the expected payoff to each agent is still strictly greater than in the absence of data because it helps rule out a low potential alternative (Nelson, 1982). When an $L$ project is revealed, there is dispersed exploration, and agents always achieve a breakthrough more efficiently than when no data is provided at all.

**Equilibrium with Data on $M$ Projects**   What is arguably more interesting, and so far understudied, is the intermediate case when a medium-value project is revealed. In this case, there exists a non-empty parameter space where data can be detrimental to exploration and social welfare due to the streetlight effect. To see this, we require the payoff from choosing M to be appealing enough relative to further exploration of other projects whose value is unknown. We restrict our attention to the following case:

**Assumption 1** ("Medium Project is Good Enough")**.**

$$M > \frac{3L + 2H}{5} \tag{3}$$

Assumption 1 ensures that selecting the medium project individually dominates searching for the high-value one. Each agent prefers to let the other agents bear the private cost of exploring, which may result in zero collective exploration and failure to achieve a breakthrough. More formally:

**Proposition 3** ("No Exploration With Data on Medium Project"). *If the underlying value of a project is revealed to be medium, then under Assumption 1 there exists a Nash equilibrium that involves all agents choosing the medium value project. Under this equilibrium, each agent's payoff is $2M$.*

*Proof.* Conditional on every other agent selecting $M$, the payoff to choosing $M$ is higher than deviating and exploring an alternative project if $2M > \frac{3}{4}(L + \max\{M, (2L+H)/3\}) + \frac{1}{4}(2H)$, which follows from Assumption 1. ∎

This means that if the value of M is tempting enough, there is always a Nash equilibrium in which zero collective exploration occurs due to the prevalence of free-riding. This result is a consequence of the informational externality that arises in our setting because agents can learn from the experimentation of others. However, the following proposition addresses the conditions under which a coordination equilibrium might still emerge:

**Proposition 4** ("Coordination May Be Possible With Data on Medium Project"). *If the underlying value of a project is revealed to be medium, then under Assumption 1 and only if*

$$M \leq \frac{H + L}{2}, \tag{4}$$

*there also exists a subgame perfect Nash equilibrium which involves agents exploring in the first period and discovering the breakthrough.*

*Proof.* A second Nash equilibrium exists where $n = 3$ agents explore and uncover $H$ by exclusion. To see this, observe that if $n = 2$ agents chose to explore in the first period, an additional agent also chooses to explore if the losses from exploring in the first period are less than the gains from exploring in the second period. This is true if $M - \frac{3}{4}L - \frac{1}{4}H \leq \frac{1}{2}(H - \max\{M, (H+L)/2\})$, which is equivalent to the condition in the proposition.

To determine the subgame perfect equilibrium, we solve by backward induction. The first $n-3$ agents to move will choose the revealed medium project in the first period. The last three agents to choose will explore new projects in the first period. Since three projects have been uncovered, the high-value project becomes known. All agents select the high-value project in the second period. ∎

For parameters that satisfy both Assumption 1 and the condition in equation (4), there exists another equilibrium solution that permits exploration and, ultimately, the discovery of a breakthrough. Under this equilibrium, some of the agents continue to select the medium project and earn $M + H$, while others are incentivized to bear the cost of exploring and earn $\frac{3}{4}L + \frac{5}{4}H$. However, this solution requires coordination among the $n$ agents: to the extent that agents fail to coordinate despite the sequential nature of choice, they will default to the original equilibrium despite not being subgame perfect. The following proposition summarizes this discussion:

**Proposition 5** ("Breakthrough With Data on Medium Project"). *If the subgame perfect coordination equilibrium survives, the following equation holds:* $P(H|M) = P(H|i) = 1$ *where* $i \in \{\emptyset, L, H\}$. *Otherwise, the following strict inequality holds:* $P(H|M) = 0 < P(H|i) = 1$.

**Welfare Conditions with Data on $M$ Projects**   Taken together, Propositions 3 and 4 shed light on the conditions that give rise to the streetlight effect. But considering how the data removes uncertainty and reveals an option that is tempting enough to forfeit exploration, is clustering around the $M$ project necessarily bad? Perhaps contrary to what one would intuitively expect, we can show that whenever $M$ is not too large relative to $L$ and $H$, then individual payoffs are lower than payoffs with no data. More formally, we introduce:

**Proposition 6** ("Individual Payoffs Under the Streetlight"). *Let $\pi_M$ be the individual payoff from the equilibrium wherein all agents select the medium value project. If*

$$M < \frac{L + 2H}{3}, \tag{5}$$

*the following strict payoff ranking holds:* $\pi_M < \pi_\emptyset < \pi_L < \pi_H$.

*Proof.* In Equation (2) of Proposition 2, we established that $\pi_\emptyset < \pi_L < \pi_H$. We only need to show that the expected individual payoff without data dominates the payoff whenever all agents select the medium-value project. This is true if $\frac{3}{5}L + \frac{1}{5}M + \frac{6}{5}H > 2M$, which is equivalent to the condition in the proposition. ∎

According to Proposition 6, the availability of data on $M$ can actually be counter-productive: payoffs are dominated by data on $L$, data on $H$, and even no data at all. As discussed above, this occurs because selecting $M$ in period 1 precludes the possibility of generating new data that could lead to a higher payoff in period 2. Despite being a payoff-maximizing choice at the individual level, the group as a whole forfeits data that would make everyone better off in the second period.

**Simple Example**   Table 1 summarizes how our results map into the parameter space of payoffs. For example, suppose we have the combination of parameters $(L, M, H) = (1, 6, 10)$. This combination satisfies both Assumption 1 and the condition in equation (5), but not the condition in equation (4). This implies the only equilibrium upon the revelation of $M$ involves all agents selecting the medium value project. In this case, $\pi_M = 12 < \pi_\emptyset = 13.8 < \pi_L = 14.5 < \pi_H = 20$. By contrast, consider now the combination of parameters $(L, M, H) = (1, 5, 10)$. While this combination also satisfies Assumption 1 and the condition in equation (5), it now does satisfy the condition in equation (4). This implies there are now multiple Nash equilibria. In one equilibrium, all agents select the medium-value

Table 1: Mapping the parameter space of our simple framework.

| Lower Limit of $M$ | Upper Limit of $M$ | Description |
|---|---|---|
| $L$ | $\frac{3L+2H}{5}$ | The value of $M$ is too low to deter exploration |
| $\frac{3L+2H}{5}$ | $\frac{H+L}{2}$ | There exist two possible equilibria:<br>1) Agents coordinate and explore<br>2) Agents do not explore<br>Only equilibrium 1) is subgame perfect. |
| $\frac{H+L}{2}$ | $\frac{L+2H}{3}$ | The only equilibrium involves no exploration.<br>The equilibrium is subgame perfect. |
| $\frac{L+2H}{3}$ | $H$ | The value of $M$ is high enough to increase welfare |

project. In this case, $\pi_M = 10 < \pi_\emptyset = 13.6 < \pi_L = 14.25 < \pi_H = 20$. However, in the second equilibrium, which is subgame perfect, individuals coordinate and explore new projects. In this case, $\pi_\emptyset = 13.6 < \pi_M = 13.95 < \pi_L = 14.25 < \pi_H = 20$. Hence, using this different set of parameters (and others that similarly satisfy the assumptions of Proposition 4), we will be able to see whether the streetlight effect continues to emerge even when exploration is expected to prevail based on the refinement criterion of subgame perfection.

# 3 Design

While our simple theoretical framework is able to rationalize the emergence of the streetlight effect, it is an open question as to whether it would explain the behavior of agents in practice. In particular, some agents might be risk-loving or pro-social, thereby exploring even when privately suboptimal to the benefit of all. It is also possible that individuals fail to calculate private payoffs or lack attention, thereby violating the key predictions of our model. To test whether these deviations are strong enough in practice to overturn our baseline results, we conducted an online experiment in which multiple participants had to solve an exploration task mirroring our theoretical framework.

## 3.1 Experimental Procedure and Logistics

Participants were invited to either the data or no-data condition in groups of ten. All ten participants logged into the platform remotely at a specific time. Upon arrival, participants received detailed written instructions about the experiment and watched a compulsory six-minute video reiterating the

main instructions while familiarizing them with the experimental platform.[1] Participants were then required to complete a short quiz as an attention and comprehension test. In addition, participants always have access to written instructions at any point in time and could contact an experimenter via cell phone or Zoom for assistance.

The experiment consisted of independent "rounds." Mimicking our conceptual framework, each round was composed of two periods over which player payoffs were calculated. The participants were randomly split into groups made of five players each. These groups were randomly reshuffled every five rounds played. In total, each participant played 20 rounds. At the end of the experiment, we collected some information on participants' demographic attributes and elicited their degree of risk aversion with a monetarily incentivized and upscaled variant of the Holt & Laury task (Holt and Laury, 2002). Participants were then paid their experimental earnings from one randomly selected round plus a show-up fee and the amount earned in the lottery associated with the risk attitude elicitation task.

The experiment was programmed with the open-source software oTree (Chen et al., 2016) and conducted by the Toronto Experimental Economics Laboratory (TEEL) and Vienna Center for Experimental Economics (VCEE). Participants were recruited from TEEL's subject pool using ORSEE (Greiner, 2015) and VCEE labs among undergraduate, master, and PhD students who had participated in at most five experiments. Participation was voluntary and participants could withdraw at any point during the experiment. We ran 35 sessions with 350 participants, and no participant was allowed to join more than one session. The age of participants ranged from 18 to 65 years, with an average of 23.21 and a standard deviation of 4.86. The experimental sessions took place in September 2021, March 2022, and July, August and October 2023. The experiment took around 50 minutes, with average earnings of CA\$ 25.97 and a standard deviation of CA\$ 1.57.[2] Additional time was needed to read the instructions, watch the explanatory video, and answer the attention quiz, so each session lasted about 75 minutes.

## 3.2 Task Description and Implementation

As shown in Panel A of Figure 1, the experimental environment and the layout of an individual round were designed to track our theoretical framework closely. Participants take the role of an individual

---

[1]The video shown to participants in the no-data condition is here: `https://drive.google.com/file/d/1TsGs2fLIcV6XFyMTmAwuJUDnAP-rAi31` and the one shown to participants in the condition with data is here: `https://drive.google.com/file/d/1vx0F-VG1P6kQQJanO99VKaYsfz1QoNbR`

[2]Participants in Vienna were paid in Euro, but we converted all amounts to CA\$ for our analysis.

**Panel A: User interface**



**Panel B: Examples of no-data condition and data conditions**

*(i) No-data condition*

*(ii) Low-value condition*



*(iii) Medium-value condition*

*(iv) High-value condition*



Figure 1: Experimental platform.

Note: This figure reproduces the interface seen by participants in our online experiment. Panel A shows how the experimental platform is seen by the participants in the no-data condition. In this example, Mountain 4 has been selected by some other participants, and the user has selected Mountain 5. Note that the dollar value of the gems changes in every round and it is showed on the left. Panel B exemplifies the four different conditions of the experiment. When subjects are assigned to the data condition, they see the value of the gem hidden behind one randomly chosen mountain. This could either be the medium, the high, or one of the low outcomes. The specific monetary value of the mapped mountain changes in every round and it is reported near the gem image.

engaged in a hunt for precious gems. There are $m = 5$ mountains, each hiding one type of gem that can only be uncovered by exploring the mountain. There are three types of gems of varying rarity and value hidden in the mountains: three topazes ($L$), one ruby ($M$), and one diamond ($H$). The exact values of the precious stones vary across rounds but the diamonds are always worth more than the rubies and the rubies are always worth more than the topazes. Participants are told that there are always three topazes, one ruby and one diamond, although they do not know which mountain hides which gem. The game's objective is to find the most valuable gems since their value directly translates into earnings in dollars.

In addition to specifying the values and distributions of the gems, the interface keeps track of the period ("stage"), the round and the "block" number as participants make their way through the experiment. A new block simply indicates the reshuffling of participants as new groups are being formed, which then stay together for five rounds. All five players in any given round are anonymous to each other and cannot interact or communicate directly.[3] Players select which mountain to explore sequentially, based on a random order that changes every round. A dynamic instruction element on their screen turns green and indicates that it is their turn to make a choice. None of them has any initial private information about the gems' location, which changes every round (but not between the first and second periods of the same round). While waiting for their turn, players can see which mountains are being selected by their co-players. When it is their turn, they can pick the same or different mountain as other players.

In the no-data condition, the two periods of a round proceed as follows. In period 1, all participants sequentially choose one mountain to explore, as described above. At the end of period 1, the gems hidden in the mountains selected by the participants are revealed to all players, and each player earns the value hidden in the mountain of their choice. In period 2, players can again choose any of the same five mountains according to the same sequential order. The position of gems remains the same, but this time participants will also see the gems located in the mountains explored in period 1. Therefore, each player can choose the same mountain of period 1 or switch to another one exploiting the new data generated by collective exploration choices. At the end of period 2, the gems contained by the mountains selected in period 2 are revealed, and their values are added to participants' round payoffs. Individual earnings for the round equal the sum of the value of the gems found in period 1 and period

---

[3]Participants know that their co-players change every five rounds, but they cannot know whom they were playing with each time, since players were not identified in any way. In a sense, players could only interact indirectly by choosing which option to explore. When a player selected an option, the other four group members only saw a generic "A group member chose this option," without ever identifying who made the choice. See Figure 1 for an example.

2 since payoffs are non-rival.

In the data condition, the two periods proceed exactly as in the no-data condition but one of the mountains is "mapped," i.e., the gem hidden behind one mountain is revealed to all participants at the start of each round. Panel B of Figure1 shows the different possibilities. Figure (i) is the no-data condition where all mountains are undisclosed. Figures (ii), (iii), and (iv) represent the three possibilities where the mapped mountain happens to have a low, medium, or high value (topaz, ruby, or diamond), respectively. Precisely which mountain is revealed and in what order is decided by a script employing a random sequence, as exemplified in Figure A.1. The data on the mapped mountain constitutes the only public information on gems' position that participants in the data condition know before starting exploring in period 1.

We collected data for 1400 rounds. Participants saw data on one of the low-value outcomes in 574 rounds, data on the medium outcome in 254 rounds, and data on the high-value outcome in 252 rounds. In the remaining 320 rounds, participants received no initial data on the gems' location.[4] Through this experiment, we used ten combinations of payoff parameters. For the experiment conducted in September 2021 and March 2022, we used the payoff for finding low, medium, and high-value gems to be one of these five combinations: $(L, M, H) = \{(1, 6, 10), (1, 6, 10.5), (1, 7, 12), (2, 7, 11), (3, 8, 12)\}$. These parameter combinations only satisfy the conditions in Assumption 1 and equation (5). Instead, for the experiment conducted in July, August and October 2023, we implemented payoff parameters that satisfy both Assumption 1 and the condition in equation (4): $(L, M, H) = \{(1, 5, 10), (1, 5, 10.5), (1, 6, 12), (2, 6, 11), (3, 7, 12)\}$. In practice, the main difference between these two sets is the value assigned to the medium outcome. Our theoretical framework predicts that rational participants should coordinate on an exploration equilibrium with the second set of parameters, preventing the streetlight effect from arising.

# 4 Experimental Results

## 4.1 Individual and Collective Outcomes

**Payoffs.** We begin by showing experimental results for all sets of parameters pooled together. For each round, we calculate the maximum possible payoff, that is, the value of the diamond times two, and compute average individual payoffs as a percentage of this value. This allows us to compare

---

[4]We used power calculations to determine the proportion of rounds with each different treatment. The final numbers are slightly asymmetric due to the random script we used to administer the experiment.

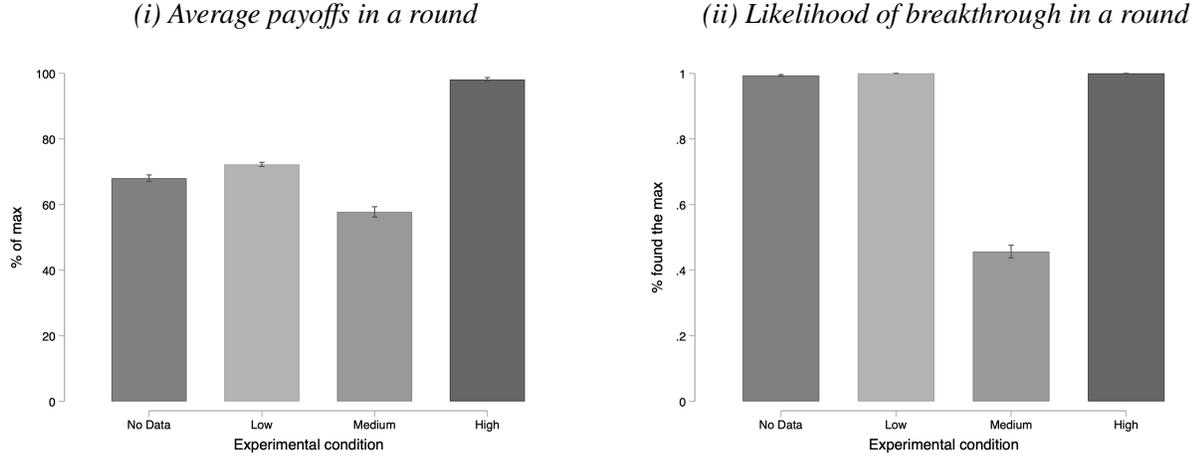| *(i) Average payoffs in a round* | *(ii) Likelihood of breakthrough in a round* |



Figure 2: Round outcomes by experimental condition.

Note: Figure (i) shows the average collective payoffs achieved in each round by experimental condition. Round payoffs are computed as a share of the maximum payoff possible in each round. Figure (ii) shows the share of rounds where the maximum was uncovered by experimental condition. Error bars represent 95% confidence intervals. See text for more details.

individual payoffs across rounds, even though the specific values of the low, medium, and high-value gems vary. We plot this average by the three data conditions and the no-data condition in Panel (i) of Figure 2. Strikingly, providing data on the medium-value project decreases average payoffs compared to all other conditions, including whenever no data is provided. We provide a precise quantification of these results by estimating the following OLS specification:

$$Payoff_{j,k} = \alpha + \beta Initial\ Data_k + \gamma \boldsymbol{X}_k + \epsilon_{j,k}, \tag{6}$$

where $Payoff_{j,k}$ is the payoff for participant $j$ in round $k$, $Data_k$ is a categorical variable encoding the type of project that was revealed at the beginning of the experiment, and $\boldsymbol{X}_k$ is a vector of controls that account for the specific payoff structure and order of appearance of the round. Standard errors are clustered at the session level. Column 1 of Table 2 presents the estimates from this regression, showing that while the average participant in the no-data condition earns about CA\$13.35, data on the medium-value mountain reduces this payoff by about 17%. Confirming the predictions of our theory, data on the high-value mountain increases payoffs by CA\$6.68, and data on the low-value mountain increases payoffs by CA\$0.89. Besides being statistically significant, these differences are also large in magnitude and thus economically meaningful.

**Breakthroughs.** Besides payoffs, the second outcome of interest is constituted by the likelihood that participants discover the high-value outcome. Panel (ii) of Figure 2 illustrates how revealing the location of the medium outcome significantly reduces the chances of a breakthrough. We quantify

Table 2: Round-level outcomes of the experiment.

| | Individual payoff (1) | I(Individual found max) (2) | I(Group found max) (3) |
|---|---|---|---|
| High | 6.682*** | 0.039** | 0.003 |
| | (0.143) | (0.011) | (0.006) |
| Low | 0.889*** | 0.037** | 0.006 |
| | (0.096) | (0.011) | (0.007) |
| Medium | -2.261*** | -0.645*** | -0.539*** |
| | (0.214) | (0.028) | (0.039) |
| Constant | 13.349*** | 0.968*** | 1.012*** |
| | (0.163) | (0.018) | (0.020) |
| Round order FE | Yes | Yes | No |
| Block order FE | Yes | Yes | Yes |
| Payoff structure FE | Yes | Yes | Yes |
| Observations | 7000 | 7000 | 1400 |

Note: $^*$ $p < 0.05$, $^{**}$ $p < 0.01$, $^{***}$ $p < 0.001$. Standard errors clustered at the session level in parentheses. Estimates from OLS models. The sample in Columns 1 and 2 is at the participant-round level (5 participants $\times$ 1400 rounds). The sample in Column 3 is at the group-round level (1400 rounds). *Individual payoff* = participant-level round payoffs in canadian dollars; *I(Individual found max)*:0/1=1 if the location of the maximum was found by the participant; *I(Group found max)*:0/1=1 if the location of the maximum was found by at least one participant in the round. The excluded category captured by the constant is the condition without data. See text for more details.

this effect using a linear probability model that mimics the specification of equation 6, but where the dependent variable is the probability of discovering the highest-value gem. Table 2 estimates that data on the medium-value project decreases the likelihood of finding the maximum of about 65% at the individual level (Column 2) and 54% at the group level (Column 3) relative to the no-data condition. Taken together, our experimental evidence aligns with the chain of predictions in Proposition 6. Innovation is highest when data directly leads to the best option, but it also increases when the available data rule out low-value alternatives. Importantly, we document that, depending on the payoff structure, the presence of more data can entail substantial societal costs.

## 4.2 Mechanisms

**Data Crowds Out Exploration.** What mechanisms explains our results? Our theoretical framework suggests that data on projects' value can lure participants into forfeiting further exploration, thus lowering long-term payoffs because of fewer discoveries. Figure 3 shows the distribution of unknown mountains selected in period 1 according to which data, if any, is present at the outset. While exploration, defined as the likelihood that an unknown mountain is chosen in period 1, is trivially low whenever the location of the maximum is known, comparing the other three conditions is informative.
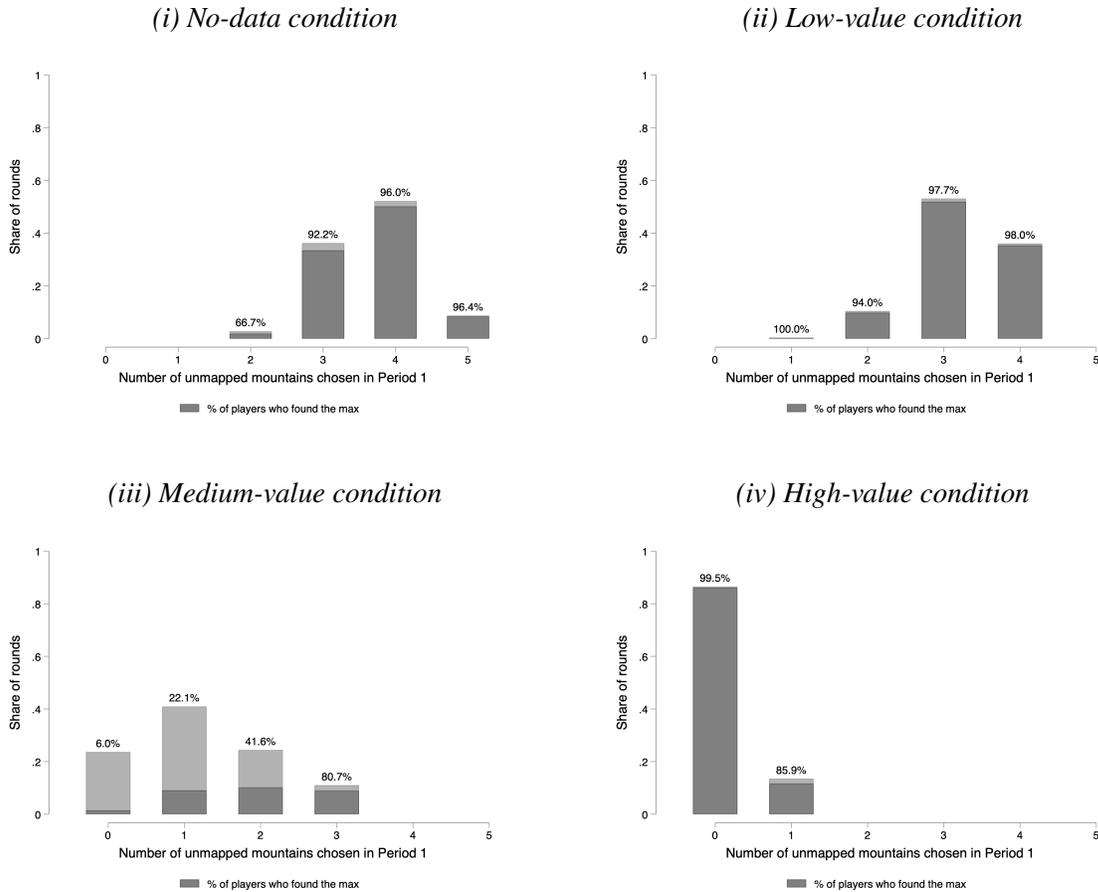
Figure 3: Distribution of the number of unknown mountains chosen in period 1 across experimental conditions, and its relationship with the discovery of the maximum.

Note: Each plot represents the empirical frequencies of rounds for each possible number of unknown options chosen in period 1, shown separately by experimental condition. In each bar, the dark grey portion shows the share of participants who found the maximum payoff in period 2. This share is also written on top of each bar to ease comparisons. See text for more details.

Receiving data on a low outcome does not reduce exploration, which remains very close to the levels of the no-data condition (Table 3, Column 1). On the contrary, when the mountain disclosed conceals the medium value outcome, collective exploration decreases by 40% relative to the no-data condition. Moreover, Panel (iii) of Figure 3 displays the positive relationship between the amount of exploration in period 1 and the share of participants that collect the maximum reward in period 2.

These exploration dynamics are also evident when the group payoffs are divided along the two periods that constitute each round (Table 3). In period 1, data on the medium outcome increases social welfare, since participants can revert to this sure option and avoid the potential failures entailed by risky experimentation. However, the situation completely reverses in period 2: the short-term gains

Table 3: Analysis of the mechanisms.

| | Exploration | Individual payoff | | I(Individual found max) | | I(Group found max) | |
|---|---|---|---|---|---|---|---|
| | Round | Period 1 | Period 2 | Period 1 | Period 2 | Period 1 | Period 2 |
| | (1) | (2) | (3) | (4) | (5) | (6) | (7) |
| High | -75.059*** | 6.499*** | 0.191* | 0.784*** | 0.037** | 0.310*** | 0.003 |
| | (1.770) | (0.097) | (0.073) | (0.008) | (0.011) | (0.015) | (0.006) |
| Low | 5.744*** | 0.664*** | 0.225** | 0.055*** | 0.036** | 0.122*** | 0.006 |
| | (1.300) | (0.074) | (0.062) | (0.007) | (0.011) | (0.020) | (0.007) |
| Medium | -34.130*** | 1.251*** | -3.511*** | -0.134*** | -0.644*** | -0.444*** | -0.539*** |
| | (2.450) | (0.122) | (0.142) | (0.007) | (0.027) | (0.029) | (0.039) |
| Constant | 83.977*** | 3.610*** | 9.717*** | 0.187*** | 0.963*** | 0.752*** | 1.012*** |
| | (2.045) | (0.104) | (0.117) | (0.008) | (0.018) | (0.018) | (0.020) |
| Round Order FE | No | Yes | Yes | Yes | Yes | No | No |
| Block order FE | Yes | Yes | Yes | Yes | Yes | Yes | Yes |
| Payoff structure FE | Yes | Yes | Yes | Yes | Yes | Yes | Yes |
| Observations | 1400 | 7000 | 7000 | 7000 | 7000 | 1400 | 1400 |

Note: $^*$ $p < 0.05$, $^{**}$ $p < 0.01$, $^{***}$ $p < 0.001$. Standard errors clustered at the session level in parentheses.

Estimates from OLS models. The sample in Column 1 is at the group-round level (1400 rounds). The sample in Columns 2, 3, 4, 5 is at the participant-period level (5 participants $\times$ 1400 periods of each type). The sample in Columns 6 and 7 is at the group-period level (1400 periods of each type). *Exploration*= share of unknown mountains explored in the round; *Individual payoff*= participant-level period payoffs in canadian dollars; *I(Individual found max)*:0/1=1 if the location of the maximum was found by the participant in the period; *I(Group found max)*:0/1=1 if the location of the maximum was found by at least one participant in the period. The excluded category captured by the constant is the condition without data. See text for more details.

from forfeiting exploration are more than offset by the cost of not uncovering the maximum. Table 3 shows that the lack of exploration early on in the game translates to a lower probability of locating the maximum, which in turn prevents its exploitation in the second period of the game. This is a direct demonstration of the streetlight effect in action: data might tilt the balance between exploration and exploitation and hurt social welfare by leaving participants stuck on a suboptimal outcome.

**Free-Riding and Coordination Failure.** In Panel (iii) of Figure 3 we showed aggregate data for all the rounds where the medium outcome was revealed. However, Proposition 4 states that a coordinated exploration is a subgame perfect equilibrium when payoffs satisfy the conditions in Assumption 1 and the condition in equation (4). Panel (i) of Figure 4 shows exploration decisions in this case, while Panel (ii) does the same for payoffs that do not satisfy the conditions for Proposition 4. We notice that there is an increase in exploration of 12.7 percentage points with the first set of payoff parameters (Appendix Table C.1). While directionally consistent with our theory, the striking result is that the general patterns are surprisingly similar between the two sets of payoffs. This finding suggests that even when collective exploration dominates the individual payoffs from free-riding, the

absence of formal coordination mechanisms in parallel experimentation settings might still prevent the emergence of an exploration equilibrium. The potential implication of this experimental result is that the streetlight effect can emerge more frequently than our theory predicts, even in contexts that, in theory, should be able to sustain collective exploration.
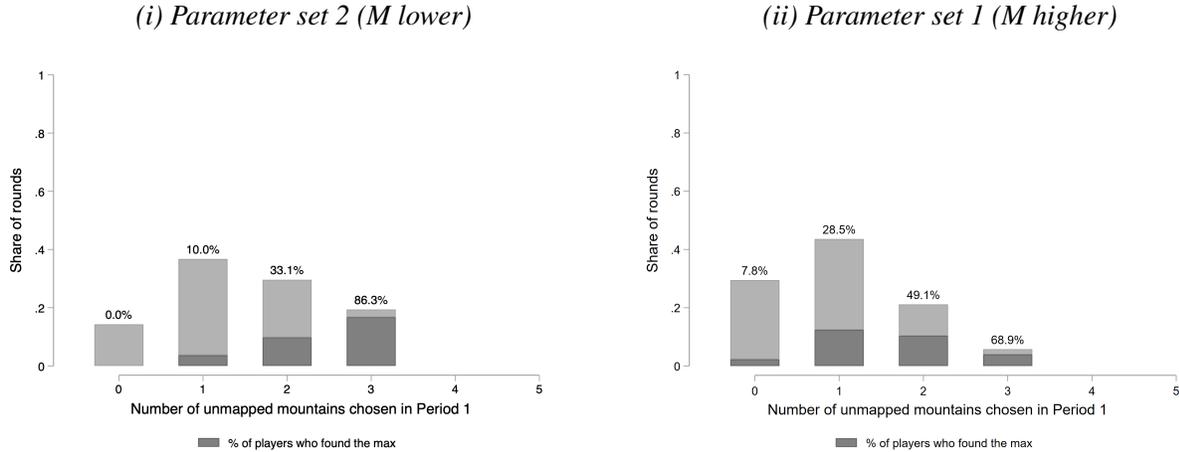


Figure 4: Distribution of the number of unknown mountains chosen in the medium-value condition in period 1, and its relationship with the discovery of the maximum.

Note: The plots represent the empirical frequencies of rounds for each possible number of unknown options chosen in period 1, shown separately for the two groups of parameters in the medium-value condition. In each bar, the dark grey portion shows the share of participants who found the maximum payoff in period 2. This share is also written on top of each bar to ease comparisons. The data underlying these figures is the same used in Panel (iii) of Figure 3, but split depending on the assumptions satisfied by each parameter set (98 rounds in Panel (i) and 156 rounds in Panel (ii)). See text for more details.

## 4.3   Exploration when Individually Suboptimal

The evidence from the experiment supports our theoretical framework, but Panel (ii) of Figure 4 shows some exploration activity even when the only Nash equilibrium involves no exploration. This finding is interesting because if enough participants deviate from the behavior predicted in Proposition 3, the streetlight effect might be prevented. In this subsection, we investigate potential explanations for why participants might choose to do so.

**Risk Aversion.**   If agents are risk neutral, as we assumed in our theoretical framework, they would always prefer the ex ante disclosed option rather than choosing an unknown option. However, it could be that risk-loving agents prefer to explore, hence explaining the empirical results. We explore this possibility using the measures of individual attitudes towards risk we collected with an incentivized variant of the Holt & Laury task at the end of the experiment (Holt and Laury, 2002). In Appendix

Table C.2, we document that risk attitudes are not associated with exploration choices when the medium outcome is known ex ante. These results suggest that participants in our experiment do not have risk preferences extreme enough to offset the negative effect of data provision, ruling out that this is what drives our findings.

**Learning.** Another explanation could be that excessive exploration is driven by an incomplete understanding of various aspects of the experiment. If this is the case, then we would also expect the number of deviations to decline over time, as participants repeatedly play variants of the game with different payoffs and mountain locations. Appendix Figure C.1 shows this is not the case. Recall that every five rounds, participants are randomly reshuffled and new groups are formed, and this procedure is repeated four times. Our results on payoffs and discovery hold for each of the four "blocks" of five rounds each. Participants behaved consistently and replicated our main results as they kept playing, without changing their propensity to select the disclosed options over time (Appendix Figures C.2 and C.3). This rules out the possibility that our results are due to limited familiarity with the experimental setup or that they would vanish as participants learn the game's dynamics.

**Correlates of Exploration Choices when Medium-Value Project is Revealed.** Appendix Table C.3 explores which individual characteristics correlate with the decision to forfeit the medium outcome in period 1. After controlling for round order, block order, and the specific payoff structure of the round in question, we find that difficulty in understanding the instructions (as proxied by being an English native speaker or giving incorrect answers to the attention quizzes) does not seem to play a major role in our setting.[5] Neither the round number nor the order in which the participant chose within the round is associated with the decision to select an unknown option.

**Qualitative Evidence.** Finally, we turn to a qualitative analysis to understand why participants' decision-making might deviate from our predictions. At the end of the experiment, we asked participants to briefly describe their thought process and how they thought the other players were making their choices. Participants' responses suggested that they grasped the game's dynamics, and the overwhelming majority described following the profit-maximizing reasoning that underlies our conceptual framework.[6] However, a handful of participants declared to follow alternative decision rules, such as

---

[5]We do find a weak correlation between exploration choices and the number of incorrect answers to the attention quizzes that followed the instructions. However, the result goes away if we add controls for the gender and age of the participant, suggesting that it is not a robust pattern. All other results are unchanged and remain insignificant once adding controls for age and gender of participants (unreported results).

[6]Notably, a few players understood the public good nature of their decisions and tried to escape the free-riding dynamics arising when the medium outcome was revealed. In this case, they realized that unless other players adopted a similar

choosing randomly, sometimes exploring an unknown option out of boredom, or taking occasional gambles. This suggests that the excessive exploration we observe in Figure 4 is due to participants' idiosyncratic preferences, but that in practice it is not strong enough to avoid the streetlight effect.

# 5 Empirical Application: The Genetic Roots of Human Diseases

The preceding sections have formalized and experimentally tested how the streetlight effect of data can emerge among agents engaged in search tasks. In what follows, we provide an empirical application that highlights how our frameworks can help rationalize the real-world patterns observed in the context of genetic research on human diseases. While our goal is not to suggest that the streetlight effect is the only reason underlying limited exploration in medical research, we investigate if the dynamics suggested by our theory are consistent with observational data from a high-stakes innovation context.

## 5.1 Setting

Knowledge about the genetic roots of human diseases dramatically increases the odds of drug development (Nelson et al., 2015). For this reason, scientists search for the genes harboring mutations that cause human diseases and that can serve as drug targets.[7] This task involves sifting through the over 19,000 protein-coding human genes, allocating experimental efforts either to the exploitation of past discoveries or the exploration of new genes. Despite individual incentives to establish priority in new areas (Bobtcheff et al., 2017; Hill and Stein, 2024), several scholars have noticed a very limited exploration of the genetic space (Edwards et al., 2011). Most research on human genes still concentrates on approximately 10% of the human genome (Gates et al., 2021). This is all the more puzzling in light of diffused awareness that potentially better drug targets exist among less explored genes (Stoeger et al., 2018). Looking for an answer, some commentators have suggested a theory akin to the streetlight effect we propose: data from past experiments on promising by ultimately futile genes might lead scientists astray (Haynes et al., 2018). In what follows, we provide an empirical investigation of this hypothesis. Our core proposition will be that disease areas where early results point to genes of moderate promise will likely see limited diversity in follow-on research and longer delays in discovering the true genetic drivers. In contrast, when early efforts yield findings of low value, higher exploration levels will help

---

strategy, the best option would be to stick to the medium-value mountain. In the words of one participant: *"If the revealed gem was a ruby, I would consider other players' choices (if they were to choose another mountain I might consider also choosing another mountain), but in most cases choosing the ruby twice gives a higher payoff."* This type of strategic behavior is consistent with participants carefully monitoring other participants whenever informational externalities are present, and ignoring their co-players when no strategic links exist (Hoelzemann and Klein, 2021).

[7]For more details about the scientific basis of this endeavor, please refer to Appendix B.

find the optimal gene-disease associations sooner.

## 5.2   Data

**DisGeNET Database.**   We assemble a dataset describing search patterns in genetics for the period of 1980-2019. The information is sourced from DisGeNET (v7.0), a complete repository of scientific publications linking human diseases to their genetic causes (Piñero et al., 2020; Tranchero, 2024). We focus on articles that explore the association between a protein-coding gene and a disease, syndrome, or abnormality with clear health implications.[8]   For each disease, we record the number of novel genetic candidates found each year as well as the yearly number of publications for previously known associations. To avoid including diseases without a genetic origin, we focus on diseases that received at least 25 publications in our sampling period.[9]   The final cross-sectional dataset includes the search and discovery histories of 4,369 diseases over a 40-year period.

**Measuring the Value of Genetic Discoveries.**   Mirroring our theoretical setup, we need a measure to classify genes as bearing $L$, $M$, or $H$ payoffs in the search for a drug target for a genetic disease. We do so by exploiting the association score assigned by DisGeNET to each gene-disease pair. The score ranges between 0 and 1, with higher values indicating greater evidence of a scientifically valuable association. The score considers the number of publications supporting the association, as well as the type and reliability of the source, with the greatest weight given to curated sources. For interpretability, we measure the scientific value of any given gene-disease association by recording its percentile in the distribution of DisGeNET scores. We then consider any score below the 60[th] percentile as a low payoff, between the 60[th] and 90[th] percentile as a medium payoff, and above the 90[th] percentile as a high payoff.[10]   Note here that only very high scores are valuable since they indicate true genetic drivers with scientific import, while middling scores are ultimately of scarce use from a therapeutic perspective. As such, the primary goal is to discover genes with a high scientific value for any given disease.[11]

---

[8]We present additional details and descriptive statistics about the DisGeNET data in Appendix B.

[9]Our results do not change if we change this threshold, see Appendix Table C.4. Furthermore, the pattern of publications is very skewed and a few diseases garnered over 5,000 publications during our sample period. Our main analysis excludes diseases in the 99[th] percentile of publications received, but we show in Appendix Table C.5 that our results are unchanged when we include such outlier diseases.

[10]In the Appendix Tables C.6 and C.7, we show that our results are robust to adopting alternative threshold values.

[11]The setting of scientific research in genetics reflects well the assumptions of non-rivalry in payoffs that we adopted in our theoretical framework. Despite the common description of scientific search as a winner-takes-all contest, recent empirical evidence finds fairly small penalties for losing a scientific race even for very narrow (and substitutable) contributions (Hill and Stein, 2024). Scientific breakthroughs usually open up avenues for follow-up work by other researchers, suggesting that their rewards are largely not rival.

**Maximum Gene Score Found During Early Search.** Our objective is to assess empirically how the presence of data on the scientific value of specific genes affects subsequent exploration patterns. Once a genetic candidate for a given disease is discovered, scientists can either concentrate on exploiting it or search for new, as-of-yet-undiscovered genetic candidates. We operationalize this idea by looking at whether the early search within a disease yielded the discovery of low, medium, or high-value genes. We define an early search window for each disease as all years before the first 10% of publications linked to the disease are reached.[12] The data available to the scientists is the value of gene-disease associations discovered during this window, as proxied by the DisGeNET score of the highest-scoring association found. We then examine how such early discovery affects future search in the disease.

**Dependent Variables.** The main dependent variable is the number of years required to achieve a breakthrough in the disease, defined as the discovery of a genetic association with a high DisGeNET score. The delay is then calculated as the years elapsed since 1980, which is the first year of our sample period.[13] To investigate the mechanisms, we look at the number of new gene candidates explored following the early search window. We account for the fact that some diseases receive a much higher amount of search efforts dividing the number of new genes explored by the number of total publications.[14] We present results for the number of new genes explored in all the years after the early search period, but results remain the same if we use a variety of alternative windows (Appendix Table C.10). Our regressions include controls for disease class (taking into account characteristics of the disease, such as its congenital or acquired origin), the total number of publications received by the disease, and the final year of the exploration window to rule out temporal variation in search activity. Estimates report standard errors clustered at the disease-class level to keep into account the correlation between similar diseases.[15]

## 5.3 Results

**Discovery of New Genetic Targets.** We begin by comparing the history of two genetic diseases, Tangier disease and Gardner syndrome (Appendix B.3). Tangier disease saw the immediate uncovering of a promising (but not optimal) genetic target in 1982. However, this discovery had the unintended

---

[12]Our results are largely consistent if we employ alternative thresholds (Appendix Table C.8) or if we use a fixed window of years to define the early search period (Appendix Table C.9).

[13]Results are unchanged if we use an alternative definition of delay that uses the years elapsed since the end of the early exploration period (Appendix Table C.11). Likewise, we obtain similar results if we look at whether a disease ever finds a breakthrough, as shown in Panel (ii) of Figure 6.

[14]In practice, this dependent variable captures changes in the average number of new genes explored for each disease.

[15]For each disease, DisGeNET indicates the disease class(es) from the Medical Subject Headings vocabulary (MeSH). Our data include 474 unique combinations of disease classes.

consequence of reducing the exploration of new genes. The result was that the true gene harboring the causal mutation was discovered only much later, in 1999. Instead, research on Gardner syndrome was initially less successful, but this led to a prolonged exploration period that culminated in the discovery of a genetic breakthrough eight years earlier than Tangier disease (Appendix Figure B.1). The dynamics highlighted by this case study apply more generally to the diseases in our sample. Diseases that initially found moderately valuable genetic targets saw slower progress toward the identification of high-value ones. Strikingly, diseases associated with low-value genes early on were 12 percentage points more likely ever to find a valuable target, taking on average 2 years less to do so (Figure 6).

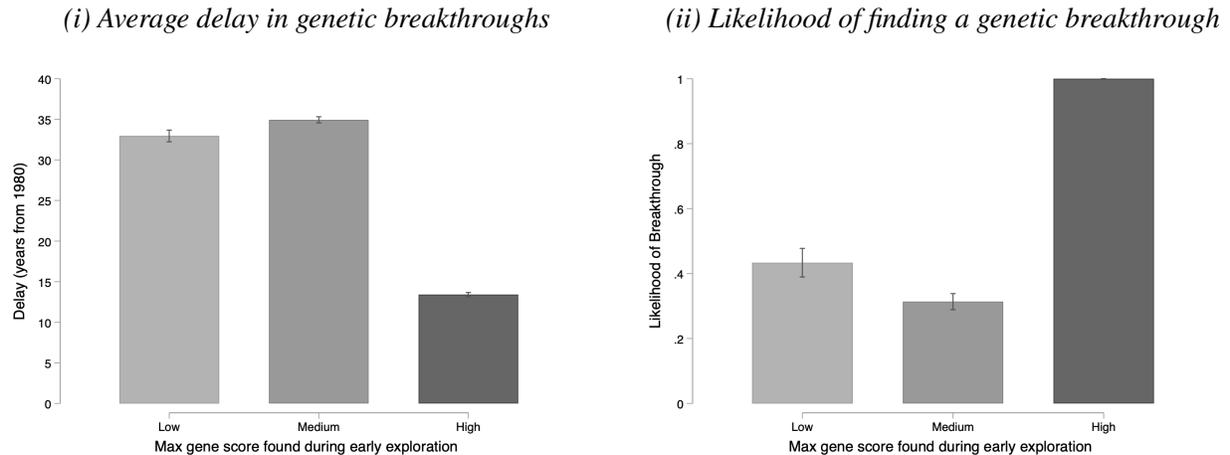*(i) Average delay in genetic breakthroughs*      *(ii) Likelihood of finding a genetic breakthrough*



Figure 6: Impact of early discoveries on search for the genetic origins of diseases.

Note: For each human disease, we compute the highest DisGeNET score identified in genetic publications linked to the disease during the early search phase (defined as the first 10% of publications on the disease). We classify maximum scores below the 60[th] percentile as a "low" gene discovery, scores between the 60[th] and 90[th] percentile as a "medium" gene discovery, and scores above the 90[th] percentile as a "high" (or breakthrough) gene discovery. Panel (i) displays the mean number of years it takes to discover the first genetic breakthrough depending on initial data. Panel (ii) displays the share of diseases that ever have a genetic breakthrough depending on initial data. Error bars represent 95% confidence intervals. See text for more details.

We explore the statistical significance of these patterns by estimating the following cross-sectional specification with OLS:

$$Delay_i = \alpha + \beta Early\ Genetic\ Data_i + \gamma \boldsymbol{X}_i + \epsilon_i, \tag{7}$$

where $Delay_i$ is the number of years from 1980 until the discovery of a genetic association with a high DisGeNET score for disease $i$, $Early\ Genetic\ Data_i$ is a categorical variable encoding the highest score found in the early search window, and $\boldsymbol{X}_{i,t}$ is a vector of controls that include the year when the early search window ends, the number of papers published about the disease, and fixed effects for disease class. The results are reported in Panel A of Table 4: the early discovery of a medium-value genetic target leads to substantially slower research progress. The estimates reveal that data on

promising but not optimal targets increase the delay by 1.6 years, or 7.2% relative to the sample mean. Note also that this delay happens despite larger search efforts following a medium finding relative to a low-value one (Appendix Figure C.4).[16]

What could explain this counterintuitive finding? Our theoretical framework offers a potential answer: the early discovery of a promising but suboptimal genetic target might have crowded out the exploration of other genes, thus resulting in slower progress despite a quantitatively larger search effort. Panel B of Table 4 shows evidence consistent with this explanation. Using the same specification in equation 7, we find that the number of new genes explored went down for diseases where a medium-value gene was found. Compared to the sample mean, the estimate implies a reduction of 15.5% in the exploration of new genes. Such a decrease is almost as large as when a breakthrough was achieved, helping to explain the subsequent delay in making high-value discoveries. While not causal in nature, this analysis provides suggestive evidence that the limited genetic exploration can be traced to dynamics consistent with our characterization of the streetlight effect (Haynes et al., 2018).[17]

**Dynamic Estimates.** The analyses above are based on cross-sectional estimates at the disease level. While proving a causal nexus between genetic data and scientific progress is beyond the scope of this paper, the richness of our dataset allows us to investigate further how exploration dynamically changes after the discovery of a medium-value opportunity to rule out key alternate explanations. If our theory is correct and the medium-value finding crowds out data generation, we should expect to observe a decrease in the proportion of new research efforts devoted to exploring new genes in the years immediately after the discovery. In practice, we test this idea by estimating the following event study specification:

$$Exploration_{i,t} = \alpha + \sum_{z} \beta_t Medium\ Gene_i \times 1(z)\ +\ \gamma \boldsymbol{X}_{i,t} + \epsilon_{i,t}, \tag{8}$$

where $Exploration_{i,t}$ is the number of new genes explored for disease $i$ in year $t$ normalized by the number of articles published, $Medium\ Gene_i \times 1(z)$ is the number of years that have elapsed since a medium value genetic association was found for disease $i$,[18] and $\boldsymbol{X}_{i,t}$ is a vector of controls that include disease fixed effects, year fixed effects, and the number of papers published each year. Figure 7 plots the coefficients resulting from this regression and shows a substantial and persistent

---

[16]The number of publications searching for the genetic roots of a given disease can be considered as a proxy for the search efforts devoted to it. We explicitly control for it in the specification presented in Column 4 of Table 4.

[17]Appendix Figure C.5 shows a parallelism between the cumulative distribution function (CDF) of our experimental results and the pattern of genetic discovery.

[18]For the few diseases with multiple medium-value genes, we define the time lags relative to the discovery of the first one.

Table 4: Impact of early genetic data on the search for the genetic origins of human diseases

**Panel A: Delay in breakthroughs**

|  | Delay (Years From 1980) | | | |
| --- | --- | --- | --- | --- |
|  | (1) | (2) | (3) | (4) |
| Max Found: Medium | 1.946*** | 1.673*** | 1.278** | 1.611** |
|  | (0.502) | (0.459) | (0.485) | (0.529) |
| Max Found: High | -18.898*** | -18.512*** | -18.538*** | -16.494*** |
|  | (0.626) | (0.611) | (0.670) | (0.726) |
| Final Exploration Year FE | No | Yes | Yes | Yes |
| Disease Class FE | No | No | Yes | Yes |
| Count of Publications | No | No | No | Yes |
| N | 3968 | 3967 | 3737 | 3337 |

**Panel B: Diversity of follow-on research**

|  | New Genes Per Paper | | | |
| --- | --- | --- | --- | --- |
|  | (1) | (2) | (3) | (4) |
| Max Found: Medium | -0.076* | -0.088** | -0.144*** | -0.115*** |
|  | (0.029) | (0.028) | (0.024) | (0.025) |
| Max Found: High | -0.306*** | -0.313*** | -0.230*** | -0.150*** |
|  | (0.028) | (0.028) | (0.031) | (0.035) |
| Final Exploration Year FE | No | Yes | Yes | Yes |
| Disease Class FE | No | No | Yes | Yes |
| Count of Publications | No | No | No | Yes |
| N | 3968 | 3967 | 3737 | 3337 |

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. Standard errors clustered at the disease class level in parentheses. Estimates from OLS models. For each human disease, we compute the highest DisGeNET score identified in the genetic publications linked to the disease during the early search phase (defined as the first 10% of publications on the disease). We classify maximum scores below the 60[th] percentile as a "low" gene discovery, scores between the 60[th] and 90[th] percentile as a "medium" gene discovery, and scores above the 90[th] percentile as a "high" (or breakthrough) gene discovery. Panel A shows the impact of early discoveries on the delay in discovering a breakthrough for a given disease, defined as years elapsed from 1980 (the first year of our panel). Panel B shows the impact of early discoveries on the number of new genes explored for a given disease, normalized by the total number of publications in the years following the exploration window. In both cases, diseases that found only low-value genes during the early search period constitute the excluded category. See text for more details.

decrease in exploratory activity after data about a medium-value opportunity become available. There is no evidence of pre-trends that foreshadow the decrease in genetic exploration observed after the discovery, further confirming the predictions of our theoretical framework.[19]  The estimates from

---

[19]Furthermore, Appendix Figure C.6 shows the same event study following the discovery of a low or high-value genetic linkage. In both cases, the dynamics resemble what our theory would imply: an increase in exploration after a low-value

a difference-in-differences specification in Appendix Table C.12 imply that after the discovery of a medium-value gene-disease association, the exploration of new genes goes down by 19% over the sample mean.

**Robustness Tests.** Our analysis documents that diseases where data highlights promising but not optimal genetic associations are less likely to make high-value discoveries by 2019 (the last year of our data). One concern is that such diseases might not have valuable genetic targets at all, implying that our cross-sectional results could be capturing this aspect and not being the reflection of actual exploration behavior. To rule out this concern, we carry out our analyses on the subset of diseases that we observe finding a high-value genetic association by 2019. Appendix Table C.13 shows that all our results are robust. A related potential issue is that, unlike our experimental set-up, genetic research might entail ambiguity on whether a breakthrough for a disease can ever be attained. Therefore, after observing a medium-value gene, scientists might not know that a better target exists among the unexplored genes. We address this issue by restricting our attention to diseases related to conditions where a breakthrough

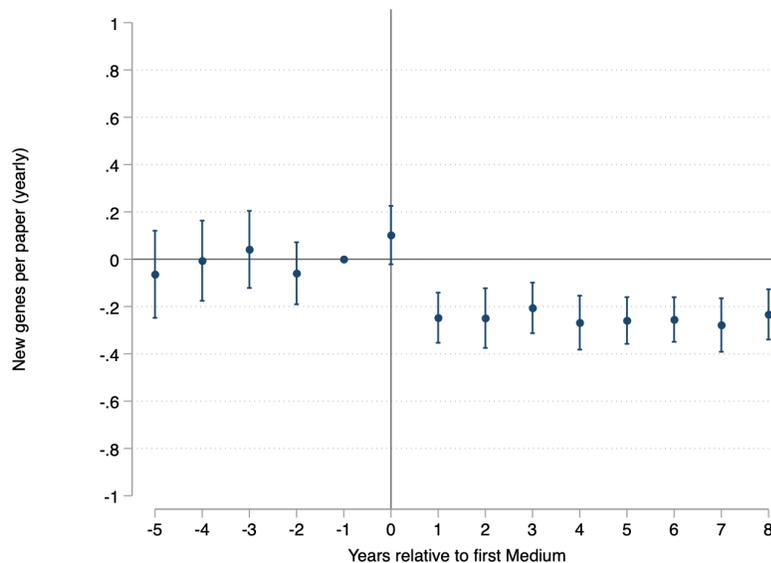discovery, and a significant decrease after the discovery of a breakthrough.



Figure 7: Dynamic effects of discovering a medium-value genetic target on subsequent exploration.

Note: For each human disease, we compute the highest DisGeNET score identified in genetic publications linked to the disease during the early search phase (defined as the first 10% of publications on the disease). We classify maximum scores between the 60th and 90th percentile as a "medium" gene discovery. This figure plots OLS estimates and 95% confidence intervals from an event study design that explores how genetic exploration in each disease evolves in the years before and after the discovery of the first medium-value genetic association. Standard errors are clustered at the disease class level. See text for more details.

gene has already been found, and for which it is thus highly likely that a breakthrough would exist as well. Appendix Figure C.7 confirms that data on a medium-value target substantially decreases exploration effort even for this subset of diseases.

# 6 Conclusion

Our paper develops a theoretical and empirical framework to understand the effects of data on exploratory search. We argued that when multiple agents need to learn about the payoffs from different projects of varying quality, providing information about "medium" quality projects can lower payoffs in the long run. Our experiment validates this prediction since we found that participants earned about 17% less when they had data about a medium outcome than when they had no data at all. Knowing which option harbors a medium outcome improves payoffs in the short run, but it reduces the likelihood that the maximum will be discovered and, therefore, lowers overall payoffs. To wit, we find that the likelihood that the optimum was collectively discovered was more than half as low when data on the medium outcome was provided. Descriptive evidence from the field of genetic research validates our idea with a consequential real-world example: early discovery of medium-value genetic targets slowed by 7.2% the discovery of the most promising genes for drug discovery.

While our work considers a simplified theoretical framework, the basic intuition could generalize to relaxing a few key assumptions. First, in our theory and experiment, we explicitly revealed a mountain of a specific value, mimicking real-world situations where agents make exploration choices after observing existing data. This approach neglects the endogeneity of data generation arising from past experimentation efforts since agents are more likely to start searching where they expect the returns to be higher and not at random. However, our main point is that even rational agents might fail to consider the informational externalities of their choices, thus contributing to creating collective equilibria of underexploration that hurt even their individual payoffs. Second, we assumed non-rivalry between participants. We posit that this assumption applies to many settings with knowledge spillovers where agents can learn from each other but are not competing directly in the product market (Bloom et al., 2013; Krieger, 2021). Our setup could also be modified to include rivalry by assuming that agents who move later receive only a fraction of the payoffs as compared to agents who choose early. If the payoff penalty applied to agents who move subsequently is relatively small, the medium value option would still be attractive for later agents, thereby maintaining our baseline predictions. This is likely the case in contexts without clear winner-takes-all dynamics, such as scientific research or technological innovation (Hill and Stein, 2024).

In sum, there is no doubt that the data revolution has the potential to dramatically lower uncertainty and boost investment in risky exploration. However, our work highlights the limits of this logic. When data points to lucrative, but ultimately less-than-ideal projects, they have the potential to do more harm than good by causing agents to herd investment activity and reduce risky exploration. Our paper thus provides support for practices such as skunkworks, where firms deliberately prevent the diffusion of initial findings of their R&D among their business units, and highlights the role of concealing intermediate information about a project unless it can be confirmed that the project represents a high-value opportunity (Boudreau and Lakhani, 2015). It also highlights the importance of policy interventions to provide "floodlights" that illuminate the entire search space, providing a common data infrastructure that prevents the streetlight effect from arising. The government can do so by providing large-scale databases, similarly to what the Landsat program did for satellite imagery (Nagaraj, 2022) or the Cancer Genome Atlas did for the genetics of cancer (Kao, 2024). As innovation and decision-making become increasingly more data-driven (Brynjolfsson and McElheran, 2016; Kim, 2023; Tranchero, 2024), more attention should be paid to strategies and policies that prevent the emergence of the streetlight effect.

# References

Ajay, A., J. McHale, and A. Oettl (2018): "Finding needles in haystacks: AI and recombinant growth," *NBER Working Paper*, 24541.

Azoulay, P., J. S. Graff Zivin, and G. Manso (2011): "Incentives and creativity: Evidence from the academic life sciences," *The RAND Journal of Economics*, 42, 527–554.

Bergeaud, A., A. Guillouzouic, E. Henry, and C. Malgouyres (2024): "From public labs to private firms: Magnitude and channels of R&D spillovers," *Banque de France, Paris School of Economics and Sciences Po*.

Bergemann, D. and A. Bonatti (2019): "Markets for information: An introduction," *Annual Review of Economics*, 11, 85–107.

Bessen, J., S. M. Impink, L. Reichensperger, and R. Seamans (2022): "The role of data for AI startup growth," *Research Policy*, 51, 104513.

Bloom, N., C. I. Jones, J. Van Reenen, and M. Webb (2020): "Are ideas getting harder to find?" *American Economic Review*, 110, 1104–1144.

Bloom, N., M. Schankerman, and J. Van Reenen (2013): "Identifying technology spillovers and product market rivalry," *Econometrica*, 81, 1347–1393.

Bobtcheff, C., J. Bolte, and T. Mariotti (2017): "Researcher's dilemma," *The Review of Economic Studies*, 84, 969–1014.

Bolton, P. and C. Harris (1999): "Strategic experimentation," *Econometrica*, 67, 349–374.

Bordalo, P., N. Gennaioli, Y. Ma, and A. Shleifer (2020): "Overreaction in macroeconomic expectations," *American Economic Review*, 110, 2748–2782.

Bordalo, P., N. Gennaioli, and A. Shleifer (2012): "Salience theory of choice under risk," *The Quarterly Journal of Economics*, 127, 1243–1285.

——— (2013): "Salience and consumer choice," *Journal of Political Economy*, 121, 803–843.

Boudreau, K. J. and K. R. Lakhani (2015): "'Open' disclosure of innovations, incentives and follow-on reuse: Theory on processes of cumulative innovation and a field experiment in computational biology," *Research Policy*, 44, 4–19.

Boyce, J. R., D. M. Bruner, and M. McKee (2016): "Strategic experimentation in the lab," *Managerial and Decision Economics*, 37, 375–391.

Brooks-Wilson, A., M. Marcil, S. M. Clee, L.-H. Zhang, K. Roomp, M. van Dam, L. Yu, C. Brewer, J. A. Collins, H. O. Molhuizen, et al. (1999): "Mutations in ABC1 in Tangier disease and familial high-density lipoprotein deficiency," *Nature Genetics*, 22, 336–345.

Brynjolfsson, E. and K. McElheran (2016): "The rapid adoption of data-driven decision-making," *American Economic Review: Papers & Proceedings*, 106, 133–39.

Camuffo, A., A. Cordova, A. Gambardella, and C. Spina (2020): "A scientific approach to entrepreneurial decision making: Evidence from a randomized control trial," *Management Science*, 66, 564–586.

Cao, R., R. Koning, and R. Nanda (2023): "Sampling bias in entrepreneurial experiments," *Management Science*.

Chen, D. L., M. Schonger, and C. Wickens (2016): "oTree—An open-source platform for laboratory, online, and field experiments," *Journal of Behavioral and Experimental Finance*, 9, 88–97.

Cockburn, I. M., R. Henderson, and S. Stern (2019): "The impact of artificial intelligence on innovation," *The Economics of Artificial Intelligence: An Agenda*, 115–152.

EDERER, F. AND G. MANSO (2013): "Is pay for performance detrimental to innovation?" *Management Science*, 59, 1496–1513.

EDWARDS, A. M., R. ISSERLIN, G. D. BADER, S. V. FRYE, T. M. WILLSON, AND F. H. YU (2011): "Too many roads not taken," *Nature*, 470, 163–165.

EWENS, M., R. NANDA, AND M. RHODES-KROPF (2018): "Cost of experimentation and the evolution of venture capital," *Journal of Financial Economics*, 128, 422–442.

FARBOODI, M. AND L. VELDKAMP (2020): "Long-run growth of financial data technology," *American Economic Review*, 110, 2485–2523.

GATES, A. J., D. M. GYSI, M. KELLIS, AND A.-L. BARABÁSI (2021): "A wealth of discovery built on the Human Genome Project—by the numbers," *Nature*, 590, 212–215.

GORDON, R. (2016): *The rise and fall of American growth: The US standard of living since the civil war*, Princeton University Press.

GREINER, B. (2015): "Subject pool recruitment procedures: organizing experiments with ORSEE," *Journal of the Economic Science Association*, 1, 114–125.

HAYNES, W. A., A. TOMCZAK, AND P. KHATRI (2018): "Gene annotation bias impedes biomedical research," *Scientific Reports*, 8, 1362.

HENRICH, J., S. J. HEINE, AND A. NORENZAYAN (2010): "Most people are not WEIRD," *Nature*, 466, 29–29.

HENRY, E., M. LOSETO, AND M. OTTAVIANI (2022): "Regulation with experimentation: Ex ante approval, ex post withdrawal, and liability," *Management Science*, 68, 5330–5347.

HILL, R. AND C. STEIN (2024): "Scooped! Estimating rewards for priority in science," *Journal of Political Economy*.

HOELZEMANN, J. AND N. KLEIN (2021): "Bandits in the lab," *Quantitative Economics*, 12, 1021–1051.

——— (2024): "Information and the bandit: The good, the bad and the ugly," *University of Vienna and University of Montreal*.

HOLT, C. A. AND S. K. LAURY (2002): "Risk aversion and incentive effects," *American Economic Review*, 92, 1644–1655.

HÖRNER, J., N. KLEIN, AND S. RADY (2022): "Overcoming free-riding in bandit games," *The Review of Economic Studies*, 89, 1948–1992.

JONES, B. F. (2009): "The burden of knowledge and the "death of the renaissance man": Is innovation getting harder?" *The Review of Economic Studies*, 76, 283–317.

JONES, C. I. AND C. TONETTI (2020): "Non-rivalry and the economics of data," *American Economic Review*, 110, 2819–2858.

KAO, J. (2024): "Charted territory: Mapping the cancer genome and R&D decisions in the pharmaceutical industry," *UCLA Anderson*.

KEHOE, A. AND V. TORVIK (2019): "Predicting Controlled Vocabulary Based on Text and Citations: Case Studies in Medical Subject Headings in MEDLINE and Patents," *University of Illinois at Urbana-Champaign*.

KELLER, G., S. RADY, AND M. CRIPPS (2005): "Strategic experimentation with exponential bandits," *Econometrica*, 73, 39–68.

KERR, W. R., R. NANDA, AND M. RHODES-KROPF (2014): "Entrepreneurship as experimentation," *Journal of Economic Perspectives*, 28, 25–48.

KIM, S. (2023): "Shortcuts to Innovation: The Use of Analogies in Knowledge Production," *Columbia Business School*.

KLEIN, N. AND S. RADY (2011): "Negatively correlated bandits," *The Review of Economic Studies*, 78, 693–732.

KRIEGER, J. L. (2021): "Trials and terminations: Learning from competitors' R&D failures," *Management Science*, 67, 5525–5548.

LEVINTHAL, D. A. (1997): "Adaptation on rugged landscapes," *Management Science*, 43, 934–950.

MANSO, G. (2011): "Motivating innovation," *The Journal of Finance*, 66, 1823–1860.

——— (2016): "Experimentation and the returns to entrepreneurship," *The Review of Financial Studies*, 29, 2319–2340.

MARCH, J. G. (1991): "Exploration and exploitation in organizational learning," *Organization Science*, 2, 71–87.

NAGARAJ, A. (2022): "The private impact of public data: Landsat satellite maps increased gold discoveries and encouraged entry," *Management Science*, 68, 564–582.

NAGARAJ, A. AND S. STERN (2020): "The economics of maps," *Journal of Economic Perspectives*, 34, 196–221.

NELSON, M. R., H. TIPNEY, J. L. PAINTER, J. SHEN, P. NICOLETTI, Y. SHEN, A. FLORATOS, P. C. SHAM, M. J. LI, J. WANG, ET AL. (2015): "The support of human genetic evidence for approved drug indications," *Nature Genetics*, 47, 856–860.

NELSON, R. R. (1982): "The role of knowledge in R&D efficiency," *Quarterly Journal of Economics*, 97, 453–470.

NGUYEN, D.-T., S. MATHIAS, C. BOLOGA, S. BRUNAK, N. FERNANDEZ, A. GAULTON, A. HERSEY, J. HOLMES, L. J. JENSEN, A. KARLSSON, ET AL. (2017): "Pharos: Collating protein information to shed light on the druggable genome," *Nucleic Acids Research*, 45, D995–D1002.

NISHISHO, I., Y. NAKAMURA, Y. MIYOSHI, Y. MIKI, H. ANDO, A. HORII, K. KOYAMA, J. UTSUNOMIYA, S. BABA, P. HEDGE, ET AL. (1991): "Mutations of chromosome 5q21 genes in FAP and colorectal cancer patients," *Science*, 253, 665–669.

PIÑERO, J., J. M. RAMÍREZ-ANGUITA, J. SAÜCH-PITARCH, RONZANO, ET AL. (2020): "The DisGeNET knowledge platform for disease genomics: 2019 update," *Nucleic Acids Research*, 48, D845–D855.

RAYMOND, L. (2024): "The Market Effect of Algorithms," *MIT Sloan Working Paper*.

SCHLAG, K. H. (1998): "Why imitate, and if so, how?: A boundedly rational approach to multi-armed bandits," *Journal of Economic Theory*, 78, 130–156.

STOEGER, T., M. GERLACH, R. I. MORIMOTO, AND L. A. NUNES AMARAL (2018): "Large-scale investigation of the reasons why potentially important genes are ignored," *PLoS Biology*, 16, e2006643.

TRANCHERO, M. (2024): "Data-driven search and innovation: Evidence from genome-wide association studies," *UC Berkeley*.

# The Streetlight Effect in Data-Driven Exploration

## Online Appendix: Additional Results

Johannes Hoelzemann

University of Vienna

Gustavo Manso

UC Berkeley-Haas

Abhishek Nagaraj

UC Berkeley-Haas & NBER

Matteo Tranchero

UC Berkeley-Haas

April 21, 2024

# A  Logistics of the Experiment

Figure A.1 summarizes how our experimental sessions unfolded. When participants join, they are assigned either to a data or to a no-data condition.[20] The experiment begins when a total of ten players are assigned to the same experimental set. Then, from each of these experimental sets, two groups of five people are randomly drawn to play the first five rounds (what we labeled as "block"). At the end of the block, the composition of the two groups is randomly reshuffled, and a second block of five rounds is played. This procedure is repeated a total of four times so that each player ends up playing exactly twenty rounds. The order of blocks seen by participants in different experimental sessions is random. The payoff structure changes each round according to a pre-recorded script generated stochastically so that the actual payoffs of each round appear random for the player. Similarly, the specific order in which specific gems are revealed in the treatment condition is generated by a random script before the experiment begins.
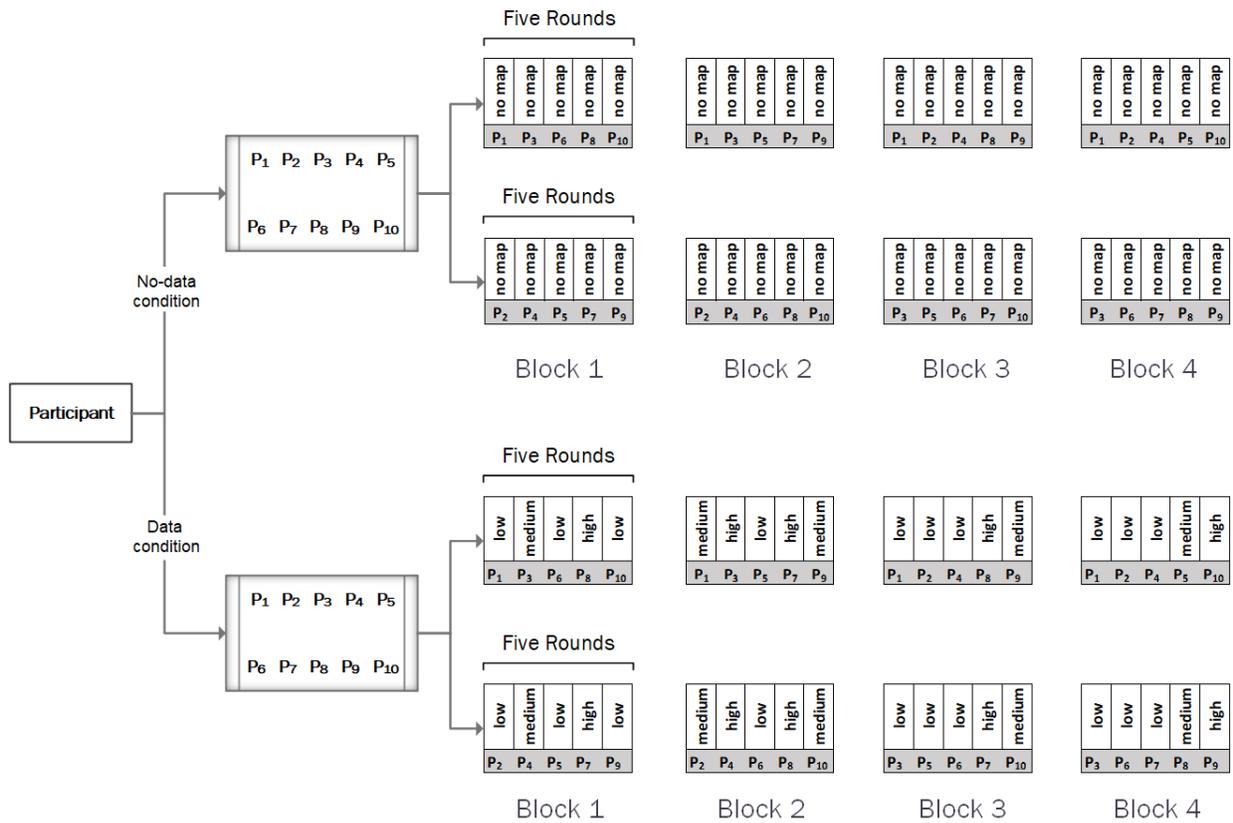


Figure A.1: Flowchart of the experimental setup.

Note: This figure provides an overview of the experiment for one actual session that took place in September 2022.

---

[20]Not in every experimental session there was a no data condition, in which case the players would be randomly split (and then reshuffled) across two distinct data conditions.

Table A.1 presents the main descriptive statistics for the experimental data, shown separately by treatment condition. Overall, the table already shows our main results in terms of payoffs and discovery while also reassuring that player characteristics were well-balanced across conditions.

Table A.1: Descriptive statistics of the experimental data.

|  | N | Mean | SD | Median | Min | Max |
|---|---|---|---|---|---|---|
| **Individual Payoff (Share)** | | | | | | |
| Low | 2870 | 72.24 | 18.63 | 62.50 | 8 | 100 |
| Median | 1270 | 57.78 | 19.01 | 58.33 | 8 | 100 |
| High | 1260 | 98.14 | 10.91 | 100.00 | 8 | 100 |
| No Data | 1600 | 68.05 | 18.93 | 62.50 | 8 | 100 |
| **I(Individual found max)** | | | | | | |
| Low | 2870 | 0.98 | 0.15 | 1.00 | 0 | 1 |
| Median | 1270 | 0.30 | 0.46 | 0.00 | 0 | 1 |
| High | 1260 | 0.98 | 0.13 | 1.00 | 0 | 1 |
| No Data | 1600 | 0.94 | 0.24 | 1.00 | 0 | 1 |
| **I(Group found max)** | | | | | | |
| Low | 574 | 1.00 | 0.00 | 1.00 | 1 | 1 |
| Median | 254 | 0.46 | 0.50 | 0.00 | 0 | 1 |
| High | 252 | 1.00 | 0.00 | 1.00 | 1 | 1 |
| No Data | 320 | 0.99 | 0.08 | 1.00 | 0 | 1 |
| **Players Age** | | | | | | |
| Low | 2870 | 22.94 | 4.43 | 22.00 | 18 | 65 |
| Median | 1270 | 22.65 | 4.13 | 22.00 | 18 | 57 |
| High | 1260 | 22.50 | 4.12 | 21.00 | 18 | 59 |
| No Data | 1600 | 24.71 | 6.18 | 23.00 | 18 | 55 |
| **Female Players** | | | | | | |
| Low | 2870 | 0.65 | 0.48 | 1.00 | 0 | 1 |
| Median | 1270 | 0.66 | 0.47 | 1.00 | 0 | 1 |
| High | 1260 | 0.65 | 0.48 | 1.00 | 0 | 1 |
| No Data | 1600 | 0.71 | 0.45 | 1.00 | 0 | 1 |

*Note*: The table presents descriptive statistics on the 7000 participants in the 1400 rounds of the experiment. *Individual Payoff Share*= individual round payoffs as a share of the maximum achievable; *I(Individual found max)*:0/1=1 if the location of the maximum was found by the participant; *I(Group found max)*:0/1=1 if the location of the maximum was found by at least one participant in the round; *Players Age*= age of the participant at the time of the experiment, in years; *Female Players*= share of participants who voluntarily reported to identify as female.

# B Searching for the Genetic Roots of Human Diseases: Additional Details

## B.1 Scientific Background

Genetics is the branch of biology that studies genes, heredity, and variation in living organisms. Genes are segments of DNA (deoxyribonucleic acid) that contain the information necessary for living organisms' development, functioning, and reproduction. In practice, each gene is a portion of DNA that contains instructions for building one or more proteins, which are the fundamental constituents of an organism. Genes often acquire mutations (or variants) in their sequence, most of which are harmless. However, some mutations can lead the gene to alter its behavior and affect phenotypic traits, sometimes with significant consequences and the emergence of severe health conditions. Discovering which mutations are responsible for specific human diseases is thus a first-order priority since genes associated with a condition can often be used as drug targets (Nelson et al., 2015). When a drug molecule binds to its genetic target, it can modify its functioning, favorably affecting the outcome of a disease. Therefore, knowing the genetic roots of diseases has important practical consequences in the design of pharmaceutical drugs.

Diseases caused by single gene mutations are called Mendelian disorders, but such diseases are typically rare. Most common human diseases have a polygenic nature, meaning they are not due to a single genetic factor but rather by many genes. This class of diseases is called complex and genetic mutations can increase the risk of presenting the condition even without being neither necessary nor sufficient. Despite often clustering in families, polygenic disorders do not have a predictable inheritance pattern because convoluted interactions between genes and environmental factors determine them. This means that scientists need to search through the over 19,000 protein-coding genes to find the mutations involved in each of the thousands of polygenic diseases (Tranchero, 2024).

Researchers have noted that even after the completion of the Human Genome Project, most scientists continue to investigate the same small number of genes (Stoeger et al., 2018). Gates et al. (2021) report that 1% of genes still receive 22% of all gene-related publications, helping to explain why current treatments exploit only around 10% of the potentially druggable targets. This situation is probably suboptimal since our chances of finding a cure for polygenic diseases would benefit from exploring a larger number of genes (Edwards et al., 2011) and several understudied genes showing high promise have been identified (Nguyen et al., 2017; Stoeger et al., 2018). Interestingly, despite much debate on this extreme concentration of attention on a small number of theoretically well-known genes, we still

lack an explanation for its drivers. Some scholars have attributed it to scientists' preference for genes with past data that permit the formulation of functional hypotheses (Haynes et al., 2018), akin to what we characterized as a streetlight effect in this paper.

## B.2 Data Description

**DisGeNET.** Our main data source is DisGeNET (v7.0), which is considered a complete repository of scientific results linking human diseases to their genetic causes (Piñero et al., 2020). This database aggregates all novel gene-disease associations studied by publications indexed in PubMed. This information is harvested from specialized sources, including curated datasets such as ClinVar, UniProt, and Orphanet.[21] In addition, DisGeNET complements these data with information extracted from the scientific literature indexed in PubMed using text-mining approaches. Our resulting data are at the gene-disease-paper level, because for each association we observe both the publication that introduced it and the list of all follow-up articles that investigated it.

**Genes.** Each gene in the database is identified by a unique identifier derived from Entrez Gene, a gene-centric database curated by the National Center for Biotechnology Information (NCBI). Entrez Gene provides tracked, unique gene identifiers that are integer and species-specific. In other words, the integer assigned to a given human gene differs from that of the homolog gene in any other species. DisGeNET compiles the Entrez Gene ID of each gene studied by papers in PubMed. We then limit our sample to protein-coding genes given their prominence in the drug discovery process (Nelson et al., 2015).

**Diseases.** Disease entries in DisGeNET are annotated using vocabulary from the Unified Medical Language System (UMLS), a set of crosswalks that bring together many health and biomedical vocabularies and standards to enable interoperability between databases. DisGeNET compiles the UMLS ID of each disease studied by papers in PubMed. Since we focus on human diseases, we keep any entries that map to the following UMLS semantic types: disease or syndrome; neoplastic process; acquired abnormality; anatomical abnormality; congenital abnormality; and mental or behavioral dysfunction. Using the UMLS ID, we also obtain disease relations from Kehoe and Torvik (2019), which contains all pairwise relationships in the Medical Subject Headings vocabulary (MeSH) hierarchy.

**Gene-Disease Pair Score.** DisGeNET provides a stable score for each gene-disease association it records. The score ranges from 0 to 1 and takes into account the number and type of sources supporting the association. In practice, it is a sum of the number of publications studying the

---

[21]For the complete list of sources aggregated by DisGeNET, see https://www.disgenet.org/dbinfo.

association, weighted by their level of curation and reliability. This synthetic metric reflects how well-established is a particular association based on current knowledge and provides a parsimonious way to assess the scientific value of any given gene-disease pair (Piñero et al., 2020; Tranchero, 2024). We then consider any score below the 60[th] percentile as a low payoff, between the 60[th] and 90[th] percentile as a medium payoff, and above the 90[th] percentile as a high payoff.[22]

**Descriptive Statistics.** We report the main descriptive statistics of our dataset in Table B.1. Panel A summarizes the data at the disease level. Around 58% of the 4,369 diseases in our sample achieved a breakthrough by 2019, which is the last year of our data. On average, it takes 22.3 years and the exploration of 131 genes to find a high-value genetic target for a disease. Panel B summarizes the data at the disease-year level. In any given year, the average disease receives 5.8 publications exploring its genetic roots, usually entailing the exploration of 3.3 new genetic associations.

Table B.1: Descriptive statistics of the DisGeNET database.

**Panel A: Disease level**

|  | Mean | Median | Sd | Min | Max | N |
|---|---|---|---|---|---|---|
| Max Found: Low (0/1) | 0.11 | 0.00 | 0.32 | 0 | 1 | 4369 |
| Max Found: Medium (0/1) | 0.31 | 0.00 | 0.46 | 0 | 1 | 4369 |
| Max Found: High (0/1) | 0.58 | 1.00 | 0.49 | 0 | 1 | 4369 |
| Year Reached 10% Papers | 2001 | 2001 | 5.55 | 1981 | 2017 | 4369 |
| Max Gene Score During Exploration | 84.79 | 92.00 | 24.27 | 0 | 100 | 4369 |
| Year of First Low Score | 1991 | 1991 | 7.08 | 1980 | 2016 | 1588 |
| Year of First Medium Score | 1995 | 1994 | 8.27 | 1980 | 2018 | 2208 |
| Year of First High Score | 1996 | 1995 | 8.24 | 1980 | 2019 | 3164 |
| Delay (Years since 1980) | 22.26 | 20.00 | 12.48 | 0 | 39 | 4369 |
| Total Publications | 229.84 | 72.00 | 526.05 | 23 | 5312 | 4369 |
| Total Genes Explored | 131.02 | 50.00 | 248.22 | 1 | 2555 | 4369 |
| New Genes Per Paper (Post-Exploration) | 0.69 | 0.68 | 0.44 | 0 | 6 | 4369 |

**Panel B: Disease-year level**

|  | Mean | Median | Sd | Min | Max | N |
|---|---|---|---|---|---|---|
| Maximum Gene Score In Year | 44.06 | 1.00 | 45.81 | 0 | 100 | 174760 |
| Yearly Count of Publications | 5.75 | 1.00 | 24.07 | 0 | 836 | 174760 |
| Yearly Count of Genes Explored | 3.28 | 0.00 | 11.36 | 0 | 646 | 174760 |

*Note*: Panel A presents descriptive statistics on papers that introduce new gene-disease associations after 2005. Panel B presents descriptive statistics of the panel dataset at the disease-year level that we used for the event-study analysis shown in Figure 7 and Appendix Figure C.6.

---

[22]As already noted in the main text, all our results are robust to the adoption of alternative threshold values (Appendix Tables C.6 and C.7).

## B.3    Case Study: Gardner's Syndrome and Tangier's Disease

To exemplify our empirical application to genetic research, consider the following two genetic diseases. Gardner syndrome (MeSH ID: D005736) is a rare disorder that falls under the umbrella of familial adenomatous polyposis. It is characterized by the development of numerous polyps, particularly in the colon and rectum. These polyps have the potential to become cancerous if left untreated. In addition to gastrointestinal manifestations, individuals with Gardner syndrome may exhibit extra-colonic features, such as the development of osteomas (benign bone tumors), particularly in the skull and jaw. Importantly, Gardner syndrome is associated with mutations in the APC gene. This tumor suppressor gene is responsible for regulating cell growth and preventing cells from dividing and multiplying too quickly.
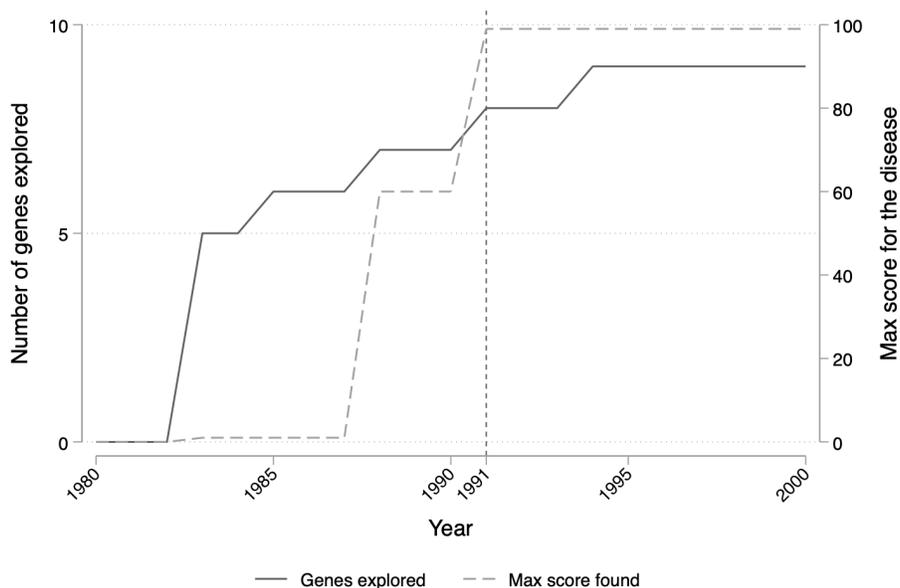
Tangier disease (MeSH ID: D013631) is a rare disorder characterized by a deficiency of high-density lipoprotein cholesterol (HDL-C) in the blood. HDL-C is responsible for transporting cholesterol away from tissues and back to the liver, playing a crucial role in cholesterol metabolism. Individuals with Tangier disease typically experience enlarged and dysfunctional tonsils that exhibit a characteristic orange discoloration. Additionally, patients may experience an increased risk of atherosclerosis and cardiovascular disease due to the decreased ability to remove cholesterol from the bloodstream. Tangier disease is an inherited genetic disease due to mutations in the ABCA1 gene. When this gene is abnormal, a problem with its instructions makes the body unable to transport lipids onto apolipoproteins, leading to a significant reduction in functional HDL-C particles.

Figure B.1 compares the history of genetic discoveries for both diseases. In the case of Gardner syndrome (Panel A), early attempts did not find any promising genes, leading to a prolonged period of exploration which culminated in the discovery of mutations in the APC gene (Nishisho et al., 1991). Instead, Tangier disease (Panel B) saw the immediate discovery of a promising association with the gene APOA1 in 1982. Such discovery stifled the exploration of new genes and led to resources being poured on a target similar to a medium-value finding in our theoretical framework (APOA1 turned out to have a DisGeNET score in the 60[th] percentile). The gene responsible for Tangier disease, ABCA1, was only discovered in 1999 by Brooks-Wilson et al. (1999).[23] This case study highlights the unfolding of the streetlight effect in a real-world example. Somewhat paradoxically, the disease for which earlier inroads were made is also the one that reached the breakthrough later. Instead, the lack of early discoveries for Gardner syndrome led to more exploration, resulting in the responsible gene being discovered 8 years earlier.

---

[23]Both these papers are very influential: Nishisho et al. (1991) and Brooks-Wilson et al. (1999) received over 2,400 and 2,100 cites in Google Scholar as of the year 2023, respectively.

**Panel A: Gardner's syndrome**
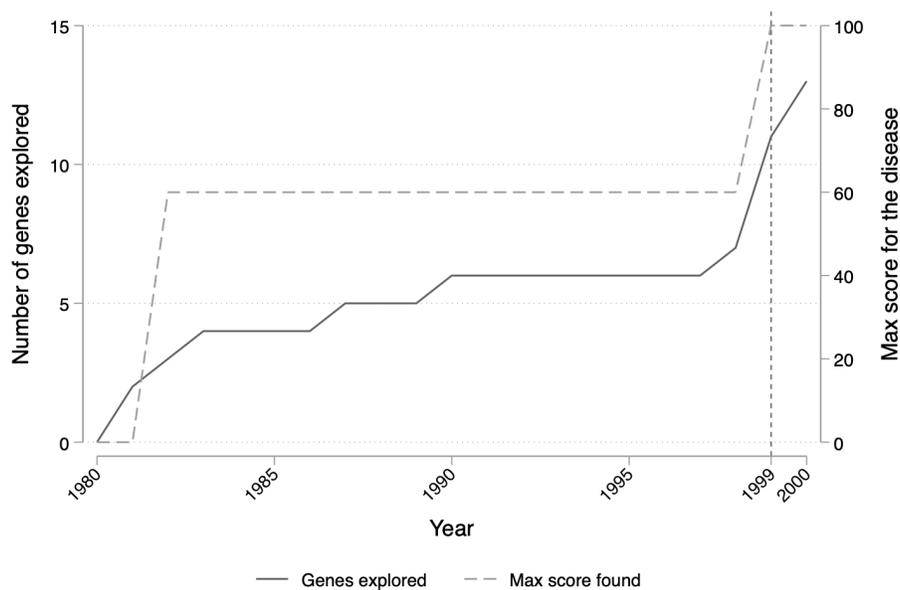


**Panel B: Tangier's disease**



Figure B.1: Two case examples in search for the genetic origins of diseases.

Note: The solid black line represents the cumulative number of gene candidates discovered for the disease up to that year. The dashed line represents the maximum DisGeNET score observed for the genes associated with that disease up to that year. Panel A presents the data for Gardner's syndrome, and the vertical line indicates the year when the association with the APC gene was discovered (DisGeNET score in the 99[th] percentile). Panel B presents the data for Tangier's disease, and the vertical line indicates the year when the association with the ABCA1 gene was discovered (DisGeNET score in the 100[th] percentile).

# C   Additional Figures and Tables

*(i) Average round payoffs by block*



*(ii) Likelihood of individual breakthrough by block*



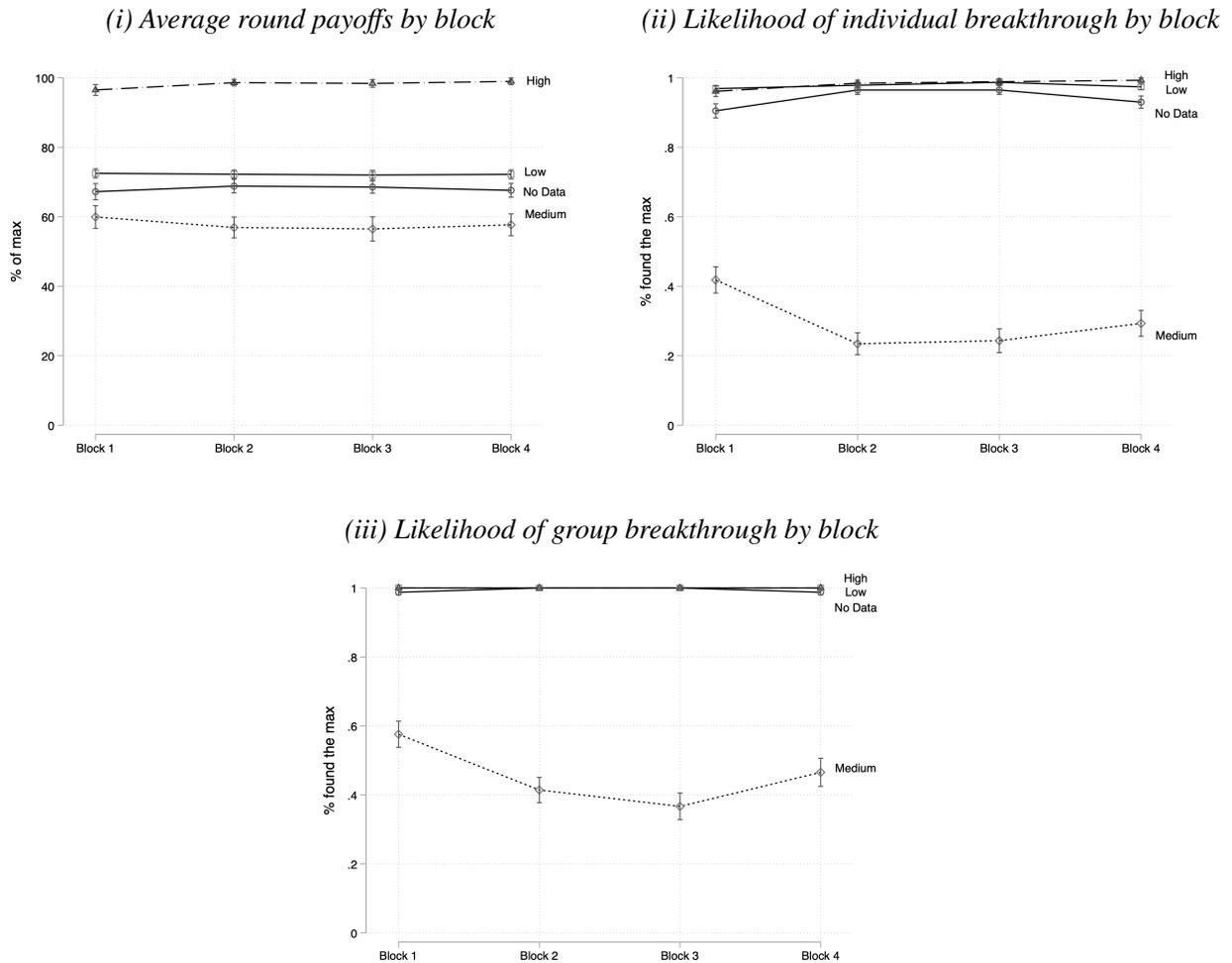*(iii) Likelihood of group breakthrough by block*



Figure C.1: Robustness of the main results over time.

Note: The figures depict the impact of data on group outcomes as the experimental session progresses. Figure (i) shows for each block of 5 rounds the average group payoffs divided by experimental condition. Payoffs are reported as a share of the maximum available in each round. Figure (ii) shows the share of participants who found the location of the maximum divided by experimental condition for each block. Figure (iii) shows the share of rounds for each block where the maximum was found divided by experimental condition.

**Panel A: Payoffs**

*(i) Payoffs in period 1*                    *(ii) Payoffs in period 2*



**Panel B: Breakthroughs**

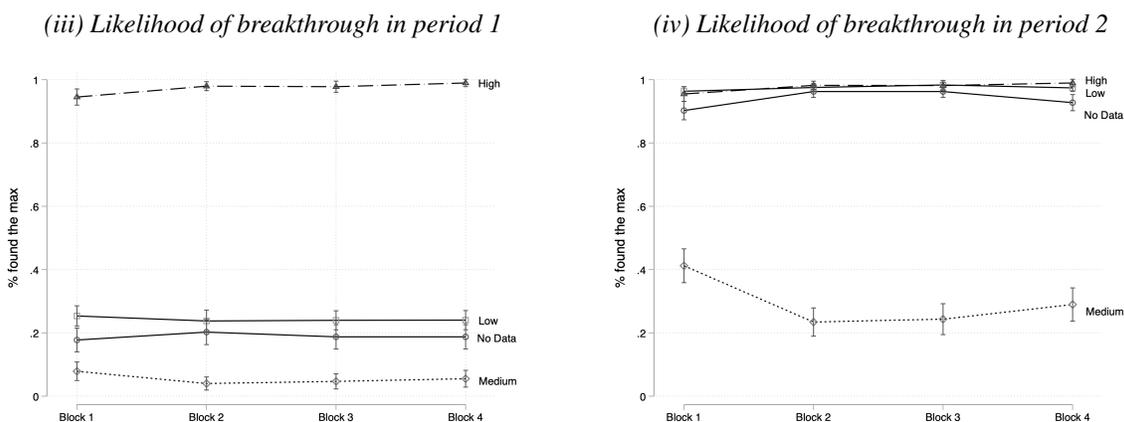*(iii) Likelihood of breakthrough in period 1*       *(iv) Likelihood of breakthrough in period 2*



Figure C.2: Outcomes over time and by period of the game.

Note: Panel A reports the experimental results on the period payoffs computed as a share of the maximum possible in each period. Figure (i) shows the average collective payoffs achieved in period 1 by experimental condition and over time. Figure (ii) shows the average collective payoffs achieved in period 2 by experimental condition and over time. Panel B reports the experimental results on the likelihood of an individual breakthrough in each round. Figure (iii) shows the share of participants that found the maximum in period 1 by experimental condition and over time. Figure (iv) shows the share of participants that found the maximum in period 2 by experimental condition and over time.
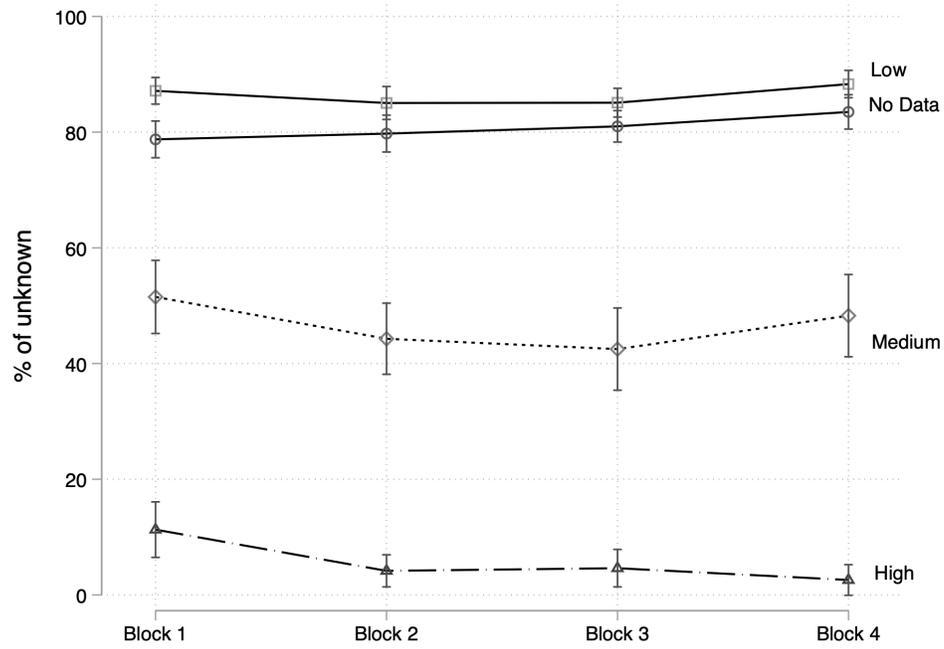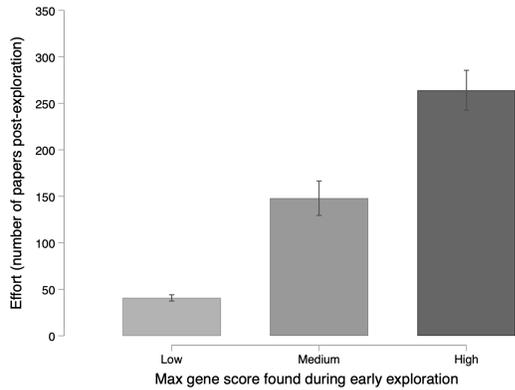
Figure C.3: Average level of exploration over time by experimental condition.

Note: The figure shows for each block of five rounds the impact of data on exploration choices divided by experimental condition. The number of mountains explored is reported as a share of the unknown mountains in each round to account for the fact that rounds without data have one more unknown option.

*(i) Average search effort for diseases*      *(ii) Average exploration of new genes for diseases*

Figure C.4: Impact of early discoveries on search effort and exploration for the genetic origins of diseases.

Note: For each human disease, we compute the highest DisGeNET score identified in genetic publications linked to the disease during the early search phase (defined as the first 10% of publications on the disease). We classify maximum scores below the 60[th] percentile as a "low" gene discovery, scores between the 60[th] and 90[th] percentile as a "medium" gene discovery, and scores above the 90[th] percentile as a "high" (or breakthrough) gene discovery. Panel (i) displays the mean number of publications about the genetic roots of a disease following the early search window. Panel (ii) displays the average number of new genes explored per paper about a disease following the early search window. Error bars represent 95% confidence intervals. See text for more details.

12

*(i) CDFs of round payoffs in the experiment*   *(ii) CDFs of delay to breakthrough in genetics*

Figure C.5: Suggestive comparison between experimental findings and patterns of genetic discovery.

Note: Panel (i) plots the cumulative density function of round payoffs by experimental condition. Panel (ii) displays the CDF of the number of years it takes to discover the first "high" score after 1980 (the first sample year) for each of the three groups. For each human disease, we compute the highest DisGeNET score identified in genetic publications linked to the disease during the early search phase (defined as the first 10% of publications on the disease). We classify maximum scores below the 60th percentile as a "low" gene discovery, scores between the 60th and 90th percentile as a "medium" gene discovery, and scores above the 90th percentile as a "high" (or breakthrough) gene discovery.

**Panel A: Discovery of a low-value genetic association**



**Panel B: Discovery of a high-value genetic association**



Figure C.6: Dynamic effects of the discovery of a low or high-value genetic association on the exploration of new genes.

Note: For each human disease, we compute the highest DisGeNET score identified in genetic publications linked to the disease during the early search phase (defined as the first 10% of publications on the disease). We classify maximum scores below the 60th percentile as a "low" gene discovery and scores above the 90th percentile as a "high" (or breakthrough) gene discovery. Panel (i) plots OLS estimates and 95% confidence intervals from an event study design that explores how genetic exploration in each disease evolves in the years before and after the discovery of the first low-value genetic association. Panel (ii) plots OLS estimates and 95% confidence intervals from an event study design that explores how genetic exploration in each disease evolves in the years before and after the discovery of the first high-value genetic association. Standard errors are clustered at the disease class level. See text for more details.

**Panel A: Keeping sibling and parent diseases**



**Panel B: Keeping only sibling diseases**



Figure C.7: Considering only diseases genetically related to a disease with a breakthrough.

Note: This figure replicates our baseline event study but only considers diseases in the sample that are genetically related to a disease with a known breakthrough (genetic discoveries with scores above the 90[th] percentile of DisGeNET score). We obtain genetic relations from the Medical Subject Headings vocabulary (MeSH). In Panel A, we keep both sibling diseases (i.e., diseases sharing the same parent MeSH code) and parent diseases (i.e., diseases one level up in the MeSH tree) of diseases with a breakthrough. In Panel B we keep only sibling diseases (i.e., diseases sharing the same parent MeSH code) of diseases with a breakthrough. This figure plots OLS estimates and 95% confidence intervals from an event study design that explores how genetic exploration in each disease evolves in the years before and after the discovery of the first medium-value genetic association. Standard errors are clustered at the disease class level. See text for more details.

Table C.1: Effects of revealing a medium-value project for different parameter values.

| | Exploration with parameter set 1 | Exploration with parameter set 2 |
|---|---|---|
| | (1) | (2) |
| | Round | Round |
| High | -75.318*** | -75.246*** |
| | (2.968) | (2.622) |
| Low | 6.473* | 4.433** |
| | (2.370) | (1.469) |
| Medium | -38.886*** | -26.168*** |
| | (4.028) | (3.123) |
| Constant | 86.290*** | 79.152*** |
| | (3.035) | (2.067) |
| Round Order FE | No | No |
| Block order FE | Yes | Yes |
| Payoff structure FE | Yes | Yes |
| Observations | 800 | 600 |

$^*\ p < 0.05,\ ^{**}\ p < 0.01,\ ^{***}\ p < 0.001.$ Standard errors clustered at the session level in parentheses

Estimates from OLS models. The sample in both columns is at the group-round level. Column 1 shows the results for rounds where the parameters used satisfy Assumption 1 and the condition in equation (5). Column 2 shows the results for rounds where the parameters used satisfy Assumption 1 and the condition in equation (4). The excluded category captured by the constant is the condition without data.

Table C.2: Risk aversion and decision not to choose the known outcome in period 1 when medium is revealed.

| | I(Exploration if M is revealed) | | |
|---|---|---|---|
| Risk aversion | -0.003 | | |
| | (0.018) | | |
| Top quartile risk aversion | | 0.000 | |
| | | (0.034) | |
| Bottom quartile risk aversion | | | -0.052 |
| | | | (0.030) |
| Constant | 0.329*** | 0.328*** | 0.343*** |
| | (0.054) | (0.051) | (0.059) |
| Round Order FE | Yes | Yes | Yes |
| Block order FE | Yes | Yes | Yes |
| Payoff structure FE | Yes | Yes | Yes |
| Observations | 1270 | 1270 | 1270 |

$^*$ $^*$ $p < 0.05$, $^{**}$ $p < 0.01$, $^{***}$ $p < 0.001$. Standard errors clustered at the session level in parentheses

Round-participant level observations, estimates from OLS models. The sample includes all the individual observations for the 254 rounds where the medium value was revealed. *I(Exploration if M is revealed)*:0/1=1 if the player did not choose the medium value in period 1. *Risk aversion* = standardized measure of individual risk aversion (Holt and Laury, 2002); *Top quartile risk aversion*:0/1=1 if the participant is in the top quartile of the risk aversion distribution in our sample; *Bottom quartile risk aversion*:0/1=1 if the participant is in the bottom quartile of the risk aversion distribution in our sample.

Table C.3: Correlates of the decision not to choose the known outcome in period 1 when medium is revealed.

| | I(Exploration if M is revealed) | | | |
|---|---|---|---|---|
| English native | 0.012 (0.038) | | | |
| Wrong quizzes | | 0.043* (0.017) | | |
| Round number | | | -0.012 (0.008) | |
| Order of choice | | | | 0.013 (0.011) |
| Constant | 0.322*** (0.067) | 0.324*** (0.054) | 0.369*** (0.055) | 0.316*** (0.056) |
| Round Order FE | Yes | Yes | Yes | No |
| Block order FE | Yes | Yes | Yes | Yes |
| Payoff structure FE | Yes | Yes | Yes | Yes |
| Observations | 1270 | 1270 | 1270 | 1270 |

$^{*} p < 0.05$, $^{**} p < 0.01$, $^{***} p < 0.001$. Standard errors clustered at the session level in parentheses

Round-player level observations, estimates from OLS models. The sample includes all the individual observations for the 254 rounds where the medium value was revealed. *I(Exploration if M is revealed)*:0/1=1 if the player did not choose the medium value in period 1. *English native*:0/1=1 if the participant is a native English speaker based on her reported nationality; *Wrong quizzes* = standardized number of wrong answers to the initial comprehension test; *Round number* = progressive order in which the rounds were played in the experimental session; *Order of choice* = random sequential order in which the player chose in that round.

Table C.4: Sensitivity to definition of marginally explored diseases

**Panel A: Delay in breakthroughs**

|  | Delay (Years From 1980) | | | |
|  | >0 Pubs | >10 Pubs | >20 Pubs | >30 Pubs |
|  | (1) | (2) | (3) | (4) |
| Max Found: Medium | 3.265*** | 2.768*** | 1.691*** | 1.362* |
|  | (0.366) | (0.392) | (0.452) | (0.602) |
| Max Found: High | -13.305*** | -14.671*** | -16.301*** | -16.931*** |
|  | (0.546) | (0.578) | (0.625) | (0.723) |
| Final Exploration Year FE | Yes | Yes | Yes | Yes |
| Disease Class FE | Yes | Yes | Yes | Yes |
| Count of Publications | Yes | Yes | Yes | Yes |
| N | 14208 | 5529 | 3828 | 2965 |

**Panel B: Diversity of follow-on research**

|  | New Genes Per Paper | | | |
|  | >0 Pubs | >10 Pubs | >20 Pubs | >30 Pubs |
|  | (1) | (2) | (3) | (4) |
| Max Found: Medium | -0.127*** | -0.129*** | -0.125*** | -0.101*** |
|  | (0.025) | (0.021) | (0.021) | (0.026) |
| Max Found: High | -0.062 | -0.125** | -0.151*** | -0.136*** |
|  | (0.042) | (0.042) | (0.035) | (0.039) |
| Final Exploration Year FE | Yes | Yes | Yes | Yes |
| Disease Class FE | Yes | Yes | Yes | Yes |
| Count of Publications | Yes | Yes | Yes | Yes |
| N | 11345 | 5529 | 3828 | 2965 |

$* p < 0.05$, $** p < 0.01$, $*** p < 0.001$. Standard errors clustered at the disease class level in parentheses. This table replicates our baseline specification (which removes diseases with less than 25 publications over the sample window) and shows robustness when we keep only diseases with nonzero publications (1), more than 10 publications (2), more than 20 publications (3), and more than 30 publications (4). For each human disease, we compute the highest DisGeNET score identified in the genetic publications linked to the disease during the early search phase (defined as the first 10% of publications on the disease). We classify maximum scores below the $60^{th}$ percentile as a "low" gene discovery, scores between the $60^{th}$ and $90^{th}$ percentile as a "medium" gene discovery, and scores above the $90^{th}$ percentile as a "high" (or breakthrough) gene discovery. Panel A shows the impact of early discoveries on the delay in discovering a breakthrough for a given disease, defined as years elapsed from 1980 (the first year of our panel). Panel B shows the impact of early discoveries on the number of new genes explored for a given disease, normalized by the total number of publications in the years following the exploration window. In both cases, diseases that found only low-value genes during the early search period constitute the excluded category. See text for more details.

Table C.5: Sensitivity to the inclusion of outlier diseases.

**Panel A: Delay in breakthroughs**

|  | Delay (Years From 1980) | | | |
|---|---|---|---|---|
|  | (1) | (2) | (3) | (4) |
| Max Found: Medium | 1.959*** | 1.701*** | 1.350** | 1.611** |
|  | (0.501) | (0.457) | (0.482) | (0.529) |
| Max Found: High | -19.043*** | -18.636*** | -18.733*** | -16.494*** |
|  | (0.647) | (0.627) | (0.694) | (0.726) |
| Final Exploration Year FE | No | Yes | Yes | Yes |
| Disease Class FE | No | No | Yes | Yes |
| Count of Publications | No | No | No | Yes |
| N | 4010 | 4009 | 3779 | 3337 |

**Panel B: Diversity of follow-on research**

|  | New Genes Per Paper | | | |
|---|---|---|---|---|
|  | (1) | (2) | (3) | (4) |
| Max Found: Medium | -0.077** | -0.089** | -0.142*** | -0.115*** |
|  | (0.029) | (0.028) | (0.024) | (0.025) |
| Max Found: High | -0.311*** | -0.316*** | -0.239*** | -0.150*** |
|  | (0.028) | (0.028) | (0.029) | (0.035) |
| Final Exploration Year FE | No | Yes | Yes | Yes |
| Disease Class FE | No | No | Yes | Yes |
| Count of Publications | No | No | No | Yes |
| N | 4010 | 4009 | 3779 | 3337 |

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. Standard errors clustered at the disease class level in parentheses. This table replicates our baseline specification including also outlier diseases (i.e., those in the top 1% by publications over the sample period). For each human disease, we compute the highest DisGeNET score identified in the genetic publications linked to the disease during the early search phase (defined as the first 10% of publications on the disease). We classify maximum scores below the 60th percentile as a "low" gene discovery, scores between the 60th and 90th percentile as a "medium" gene discovery, and scores above the 90th percentile as a "high" (or breakthrough) gene discovery. Panel A shows the impact of early discoveries on the delay in discovering a breakthrough for a given disease, defined as years elapsed from 1980 (the first year of our panel). Panel B shows the impact of early discoveries on the number of new genes explored for a given disease, normalized by the total number of publications in the years following the exploration window. In both cases, diseases that found only low-value genes during the early search period constitute the excluded category. See text for more details.

Table C.6: Alternative definitions of low and medium-value genes.

**Panel A: Delay in breakthroughs**

| | Delay (Years From 1980) | | | |
|---|---|---|---|---|
| | $50^{th}$ P | $60^{th}$ P | $70^{th}$ P | $80^{th}$ P |
| | (1) | (2) | (3) | (4) |
| Max Found: Medium | 1.924** | 1.611** | 1.611** | 1.568** |
| | (0.685) | (0.529) | (0.529) | (0.557) |
| Max Found: High | -16.044*** | -16.494*** | -16.494*** | -16.858*** |
| | (0.899) | (0.726) | (0.726) | (0.626) |
| Final Exploration Year FE | Yes | Yes | Yes | Yes |
| Disease Class FE | Yes | Yes | Yes | Yes |
| Count of Publications | Yes | Yes | Yes | Yes |
| N | 3337 | 3337 | 3337 | 3337 |

**Panel B: Diversity of follow-on research**

| | New Genes Per Paper | | | |
|---|---|---|---|---|
| | $50^{th}$ P | $60^{th}$ P | $70^{th}$ P | $80^{th}$ P |
| | (1) | (2) | (3) | (4) |
| Max Found: Medium | -0.055 | -0.115*** | -0.115*** | -0.183*** |
| | (0.031) | (0.025) | (0.025) | (0.025) |
| Max Found: High | -0.113*** | -0.150*** | -0.150*** | -0.159*** |
| | (0.030) | (0.035) | (0.035) | (0.032) |
| Final Exploration Year FE | Yes | Yes | Yes | Yes |
| Disease Class FE | Yes | Yes | Yes | Yes |
| Count of Publications | Yes | Yes | Yes | Yes |
| N | 3337 | 3337 | 3337 | 3337 |

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. Standard errors clustered at the disease class level in parentheses. This table replicates our baseline specification but varies the cutoff between a low and medium-value genetic association. In our baseline, we adopt the 60th percentile to separate medium and high scores. We test the 50th percentile (1), the baseline (2), the 70th percentile (3), and the 80th percentile (4) instead. For each regression, we hold the cutoff between a medum gene score and high gene score fixed at the 90th percentile (our baseline). Panel A shows the impact of early discoveries on the delay in discovering a breakthrough for a given disease, defined as years elapsed from 1980 (the first year of our panel). Panel B shows the impact of early discoveries on the number of new genes explored for a given disease, normalized by the total number of publications in the years following the exploration window. In both cases, diseases that found only low-value genes during the early search period constitute the excluded category. See text for more details.

Table C.7: Alternative definitions of medium and high-value genes.

**Panel A: Delay in breakthroughs**

|  | Delay (Years From 1980) | | | |
| --- | --- | --- | --- | --- |
|  | $90^{th}$ P | $95^{th}$ P | $98^{th}$ P | $99^{th}$ P |
|  | (1) | (2) | (3) | (4) |
| Max Found: Medium | 1.611** | 0.371 | 0.494 | 0.355 |
|  | (0.529) | (0.529) | (0.541) | (0.512) |
| Max Found: High | -16.494*** | -17.936*** | -18.299*** | -18.954*** |
|  | (0.726) | (0.793) | (0.803) | (0.821) |
| Final Exploration Year FE | Yes | Yes | Yes | Yes |
| Disease Class FE | Yes | Yes | Yes | Yes |
| Count of Publications | Yes | Yes | Yes | Yes |
| N | 3337 | 3337 | 3337 | 3337 |

**Panel B: Diversity of follow-on research**

|  | New Genes Per Paper | | | |
| --- | --- | --- | --- | --- |
|  | $90^{th}$ P | $95^{th}$ P | $98^{th}$ P | $99^{th}$ P |
|  | (1) | (2) | (3) | (4) |
| Max Found: Medium | -0.115*** | -0.083*** | -0.071** | -0.072** |
|  | (0.025) | (0.024) | (0.024) | (0.025) |
| Max Found: High | -0.150*** | -0.220*** | -0.249*** | -0.275*** |
|  | (0.035) | (0.039) | (0.040) | (0.039) |
| Final Exploration Year FE | Yes | Yes | Yes | Yes |
| Disease Class FE | Yes | Yes | Yes | Yes |
| Count of Publications | Yes | Yes | Yes | Yes |
| N | 3337 | 3337 | 3337 | 3337 |

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. Standard errors clustered at the disease class level in parentheses. This table replicates our baseline specification but varies the cutoff between a medium and high-value genetic association. In our baseline, we adopt the 90th percentile to separate medium and high scores. We test the baseline (1), the 95th percentile (2), the 98th percentile (3), and the 99th percentile (4) instead. For each regression, we hold the cutoff between a low gene score and medium gene score fixed at the 60th percentile (our baseline). Panel A shows the impact of early discoveries on the delay in discovering a breakthrough for a given disease, defined as years elapsed from 1980 (the first year of our panel). Panel B shows the impact of early discoveries on the number of new genes explored for a given disease, normalized by the total number of publications in the years following the exploration window. In both cases, diseases that found only low-value genes during the early search period constitute the excluded category. See text for more details.

Table C.8: Different thresholds of publication share to define the early search period.

| | Delay (Years From 1980) | | | |
|---|---|---|---|---|
| | 5% | 10% | 15% | 20% |
| | (1) | (2) | (3) | (4) |
| Max Found: Medium | 2.598*** | 1.611** | 0.661 | 1.008 |
| | (0.436) | (0.529) | (0.550) | (0.548) |
| Max Found: High | -13.657*** | -16.494*** | -18.465*** | -18.532*** |
| | (0.680) | (0.726) | (0.686) | (0.688) |
| Final Exploration Year FE | Yes | Yes | Yes | Yes |
| Disease Class FE | Yes | Yes | Yes | Yes |
| Count of Publications | Yes | Yes | Yes | Yes |
| N | 3325 | 3337 | 3369 | 3391 |

$^*$ $p < 0.05$, $^{**}$ $p < 0.01$, $^{***}$ $p < 0.001$. Standard errors clustered at the disease class level in parentheses. This table replicates our baseline specification using alternative definitions of "early search period". We test varying thresholds from 5% (1) up to 20% (4) in increments of 5%. This table replicates our baseline specification using alternative windows to define the period of early search. We report the results employing fixed windows, including all years before 1990 (1), before 1995 (2), before 2000 (3), and before 2005 (4). For each human disease, we compute the highest DisGeNET score identified in the genetic publications linked to the disease during the early search phase. We classify maximum scores below the 60[th] percentile as a "low" gene discovery, scores between the 60[th] and 90[th] percentile as a "medium" gene discovery, and scores above the 90[th] percentile as a "high" (or breakthrough) gene discovery. See text for more details.

Table C.9: Fixed windows of years to define the early search period.

| | Delay (Years From 1980) | | | |
|---|---|---|---|---|
| | <1990 | <1995 | <2000 | <2005 |
| | (1) | (2) | (3) | (4) |
| Max Found: Medium | 1.939** | 1.811** | 2.531*** | 1.560** |
| | (0.724) | (0.556) | (0.544) | (0.521) |
| Max Found: High | -18.706*** | -17.826*** | -18.596*** | -20.288*** |
| | (0.958) | (0.778) | (0.653) | (0.581) |
| Final Exploration Year FE | Yes | Yes | Yes | Yes |
| Disease Class FE | Yes | Yes | Yes | Yes |
| Count of Publications | Yes | Yes | Yes | Yes |
| N | 1192 | 2213 | 2923 | 3297 |

$^{*}$ $p < 0.05$, $^{**}$ $p < 0.01$, $^{***}$ $p < 0.001$. Standard errors clustered at the disease class level in parentheses. This table replicates our baseline specification using alternative windows to define the period of early search. We report the results employing fixed windows, including all years before 1990 (1), before 1995 (2), before 2000 (3), and before 2005 (4). For each human disease, we compute the highest DisGeNET score identified in the genetic publications linked to the disease during the early search phase. We classify maximum scores below the $60^{th}$ percentile as a "low" gene discovery, scores between the $60^{th}$ and $90^{th}$ percentile as a "medium" gene discovery, and scores above the $90^{th}$ percentile as a "high" (or breakthrough) gene discovery. See text for more details.

Table C.10: Alternative windows to examine follow-on explorative research.

| | New Genes Per Paper | | | |
| | All Years | 5 Years | 10 Years | Until H |
| | (1) | (2) | (3) | (4) |
|---|---|---|---|---|
| Max Found: Medium | -0.115*** | -0.098* | -0.105** | -0.095* |
| | (0.025) | (0.039) | (0.033) | (0.040) |
| Max Found: High | -0.150*** | -0.183*** | -0.179*** | |
| | (0.035) | (0.044) | (0.041) | |
| Disease Class FE | Yes | Yes | Yes | Yes |
| Final Exploration Year FE | Yes | Yes | Yes | Yes |
| Count of Publications | Yes | Yes | Yes | Yes |
| N | 3337 | 3305 | 3332 | 1077 |

$^*$ $p < 0.05$, $^{**}$ $p < 0.01$, $^{***}$ $p < 0.001$. Standard errors clustered at the disease class level in parentheses. This table replicates our baseline specification using alternative windows to evaluate the evolution of explorative research. We report the results from the baseline (1), the 5 subsequent years after 10% of publications is reached (2), the 10 subsequent years after 10% of publications is reached (3), and until the first high gene score is found (4). For each human disease, we compute the highest DisGeNET score identified in the genetic publications linked to the disease during the early search phase (defined as the first 10% of publications on the disease). We classify maximum scores below the 60th percentile as a "low" gene discovery, scores between the 60th and 90th percentile as a "medium" gene discovery, and scores above the 90th percentile as a "high" (or breakthrough) gene discovery. See text for more details.

Table C.11: Alternative measures of delays in breakthroughs.

| | Delay (Years From 10% Publications) | | |
|---|---|---|---|
| | (1) | (2) | (3) |
| Max Found: Medium | 4.063*** | 3.182** | 3.302** |
| | (1.003) | (0.994) | (0.994) |
| Max Found: High | -22.939*** | -22.216*** | -21.194*** |
| | (1.118) | (1.250) | (1.242) |
| Disease Class FE | No | Yes | Yes |
| Count of Publications | No | No | Yes |
| N | 3968 | 3738 | 3338 |

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. Standard errors clustered at the disease class level in parentheses. This table replicates our baseline specification using an alternative measure of delay in breakthrough discovery, here defined as years elapsed from the first 10% of publications on the disease. For each human disease, we compute the highest DisGeNET score identified in the genetic publications linked to the disease during the early search phase (defined as the first 10% of publications on the disease). We classify maximum scores below the 60th percentile as a "low" gene discovery, scores between the 60th and 90th percentile as a "medium" gene discovery, and scores above the 90th percentile as a "high" (or breakthrough) gene discovery. See text for more details.

Table C.12: Difference-in-difference estimates of the effect of early discoveries on subsequent genetic exploration.

| | New Genes Paper Paper (Yearly) | | | | | |
| | Low GDA | | Medium GDA | | High GDA | |
| | (1) | (2) | (3) | (4) | (5) | (6) |
| --- | --- | --- | --- | --- | --- | --- |
| Post (Low) | 0.252*** | 0.225*** | | | | |
| | (0.039) | (0.039) | | | | |
| Post (Med) | | | -0.139*** | -0.141*** | | |
| | | | (0.028) | (0.029) | | |
| Post (High) | | | | | -0.177*** | -0.153*** |
| | | | | | (0.019) | (0.019) |
| Disease FE | Yes | Yes | Yes | Yes | Yes | Yes |
| Year FE | Yes | Yes | Yes | Yes | Yes | Yes |
| Count of Publications | No | Yes | No | Yes | No | Yes |
| N | 88134 | 87997 | 88134 | 87997 | 88134 | 87997 |

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. Standard errors clustered at the disease class level in parentheses. For each human disease, we compute the highest DisGeNET score identified in the genetic publications linked to the disease during the early search phase (defined as the first 10% of publications on the disease). We classify maximum scores below the 60th percentile as a "low" gene discovery, scores between the 60th and 90th percentile as a "medium" gene discovery, and scores above the 90th percentile as a "high" (or breakthrough) gene discovery. This table reports OLS estimates from differences-in-differences that explore how genetic exploration in each disease evolves in the years before and after the discovery of the first low, medium, and high-value genetic association. Standard errors are clustered at the disease class level. See text for more details.

Table C.13: Considering only diseases that have a breakthrough by 2019.

**Panel A: Delay in breakthroughs**

| | Delay (Years From 1980) | | | |
| --- | --- | --- | --- | --- |
| | (1) | (2) | (3) | (4) |
| Max Found: Medium | 0.918 | 0.759 | 0.641 | 2.199*** |
| | (0.594) | (0.507) | (0.560) | (0.473) |
| Max Found: High | -11.591*** | -11.800*** | -11.832*** | -8.531*** |
| | (0.536) | (0.521) | (0.591) | (0.441) |
| Final Exploration Year FE | No | Yes | Yes | Yes |
| Disease Class FE | No | No | Yes | Yes |
| Count of Publications | No | No | No | Yes |
| N | 3053 | 3051 | 2861 | 2477 |

**Panel B: Diversity of follow-on research**

| | New Genes Per Paper | | | |
| --- | --- | --- | --- | --- |
| | (1) | (2) | (3) | (4) |
| Max Found: Medium | 0.005 | 0.004 | -0.071 | -0.067 |
| | (0.044) | (0.044) | (0.043) | (0.042) |
| Max Found: High | -0.210*** | -0.220*** | -0.181*** | -0.113** |
| | (0.037) | (0.038) | (0.031) | (0.042) |
| Final Exploration Year FE | No | Yes | Yes | Yes |
| Disease Class FE | No | No | Yes | Yes |
| Count of Publications | No | No | No | Yes |
| N | 3053 | 3051 | 2861 | 2477 |

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. Standard errors clustered at the disease class level in parentheses. This table replicates our baseline specification removing any diseases without a breakthrough (i.e. a gene with a "high" score) during the sample period. For each human disease, we compute the highest DisGeNET score identified in the genetic publications linked to the disease during the early search phase (defined as the first 10% of publications on the disease). We classify maximum scores below the 60th percentile as a "low" gene discovery, scores between the 60th and 90th percentile as a "medium" gene discovery, and scores above the 90th percentile as a "high" (or breakthrough) gene discovery. Panel A shows the impact of early discoveries on the delay in discovering a breakthrough for a given disease, defined as years elapsed from 1980 (the first year of our panel). Panel B shows the impact of early discoveries on the number of new genes explored for a given disease, normalized by the total number of publications in the years following the exploration window. In both cases, diseases that found only low-value genes during the early search period constitute the excluded category. See text for more details.