

ORGANIZING MODULAR PRODUCTION*

Niko Matouschek

Northwestern University

Michael Powell

Northwestern University

Bryony Reich

Northwestern University

November 8, 2021

Preliminary and Incomplete

Abstract

We characterize the optimal communication network in a firm with a modular production function, which we model as a network of decisions with a non-overlapping community structure. Optimal communication is characterized by two nested hierarchies that determine whom each agent receives information from and sends information to. Receiver rank depends only on the cohesion of an agent's module while sender rank also depends on the decision-specific value of adaptation. We provide conditions for the Mirroring Hypothesis to hold and show that optimal communication is consistent with a core-periphery structure in which the core is formed by the most cohesive modules.

Keywords: network design, modularity, communication, hierarchies

JEL classifications: D23, D85, L23

*We thank Heski Bar-Isaac, Hector Chade, Krishna Dasaratha, Wouter Dessein, Itay Fainmesser, Ben Golub, Matt Jackson, Willemien Kets, Andrea Prat, Jeroen Swinkels, Alireza Tahbaz-Salehi, Eduard Talamàs, and participants at various conferences and seminars for comments and suggestions. We also thank Hossein Alidaee for excellent research assistance. Matouschek: n-matouschek@kellogg.northwestern.edu; Powell: mike-powell@kellogg.northwestern.edu; Reich: bryony.reich@kellogg.northwestern.edu

1 Introduction

Over the last 60 years, the economy has experienced a remarkable shift towards modular production (Baldwin and Clark 2000). Nowadays so many products are made by assembling separately produced modules that the 21st century has been called the *Modular Age* (see, for instance, Garud, Kumaraswamy, and Langlois (2009)). The rise of modular production has the potential to change the organization of firms, the structure of industries, and the location of production. In this paper we take a first step towards exploring the economic implications of modular production by examining its impact on the internal organization of firms.

Herbert Simon anticipated the rise of modular production in 1962, when he observed that complex social, technological, and biological systems—large firms, mechanical watches, the human body—tend to be made up of communities or modules, groups of elements with stronger within than across group interactions (Simon 1962). The advantage of this modular structure, he argued, is that it allows systems to adapt to changes in the environment by making adjustments in a limited number of modules while leaving the rest of the system unchanged. The prevalence of modular structures has since been confirmed by the literature on community detection, which has documented them in a wide variety of contexts from the internet to the global air transportation network and the brain (Meunier et al. 2009, Guimera et al. 2005, and Fortunato 2010).

A few years after Simon wrote his article, IBM developed the first modular computer, the System/360. Until then computers had been tightly integrated systems of their constituent parts. A change in the processor or any other critical component required the design of an entirely new computer. This made it difficult to adopt new technologies and adapt computers to the idiosyncratic demands of different customers. The System/360 was designed to change all this. Its modular structure was a deliberate choice by IBM’s executives who tasked their engineers with developing a computer that was made up of a small number of easily assemblable and exchangeable modules. Henceforth, when a supplier developed a better disk drive, or a customer needed more storage, IBM was able to adapt quickly. Not only did this make the System/360 an enormous financial success, it also changed how computers have been built ever since (Baldwin and Clark 1997 and 2000).

The move towards modular production has not been confined to the computer industry. Over the last few decades, firms across a wide range of industries followed in IBM’s footsteps and developed products with modular production functions. Smartphones, airplanes, and electric cars are all made by assembling a limited number of modules. Even homes are now routinely assembled

from pre-made modules rather than built on site from scratch. Nor is this move towards modular production confined to physical products. Modular programming—separating the different functions of computer programs into independent and interchangeable modules—was developed in the 1960s and is now a common feature of almost all programming languages. Of course, many products exhibited some degree of modularity even before the System/360. Builders installed pre-made doors and windows long before the rise of modular home building. What is different now, though, is that many products are modular by design. They are produced entirely by assembling a limited number of modules and, in line with Herbert Simon’s observations, they are now more rule than exception.¹

The widespread adoption of modular production has the potential to change the organization of production and thus the outcomes of economic activity. In the short run, firms adapt their internal organizations to accommodate modular production. Over time, they may also change their boundaries which, in turn, can alter the structure of their industries and the location of production. To manage the System/360, for instance, IBM established a centralized office, which ensured that different modules worked together, but also delegated control over individual modules to autonomous teams. This process of decentralization continued over many years with IBM and its competitors eventually outsourcing the development and production of modules to smaller, independent, and often foreign firms (Baldwin and Clark 1997 and 2000).

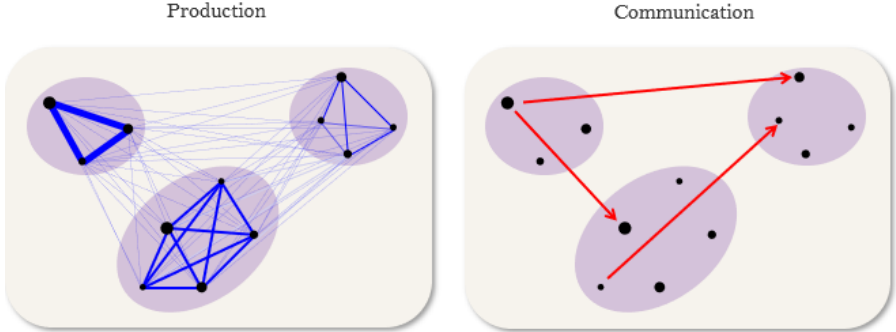


Figure 1: Left panel—The firm’s production function takes the form of a network with a non-overlapping community structure. Right panel—Given the production network, the principal designs the optimal communication network by deciding whom each agent should tell his state to, taking as given that each directed link comes at an exogenous cost.

The goal of this paper is to take a first step towards exploring the economic implications

¹See Baldwin and Clark (1997) and (2000). See also the Wikipedia entries for Modular Design, Modular Programming, and Modular Building and the references therein.

of modular production (“first step” in the economics literature; as discussed below there is an expansive existing literature in management). We focus on the internal organization of a single firm with a modular production function, which we model as a network of decisions with a non-overlapping community structure and illustrate in Figure 1. Every node represents a decision, an agent who makes the decision, and a state. The size of the node represents the importance of adapting the decision to its state, and the width of an edge between two decisions represents the importance of coordinating the two decisions. The decisions are partitioned into “modules,” groups of decisions that require more coordination with each other than with decisions in other modules and that are indicated by the shaded areas in the figure. The adjacency matrix of the production network, therefore, takes the form of a block matrix. This structure approximates the interactions between decisions in modular products, such as the laptop computer illustrated in Figure 2

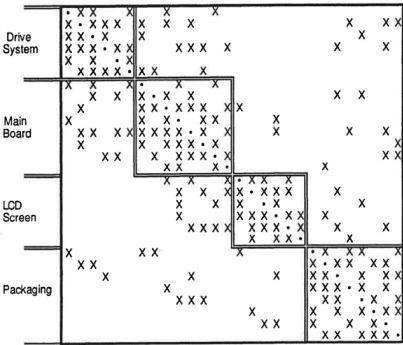


Figure 2: The “Design Structure Matrix” of a laptop computer in which each row and column corresponds to a task in producing the computer and an “x” entry indicates a strong need for coordination (reproduction of Figure 2.3 in McCord and Eppinger (1993)).

In the model, the only impediment to efficient production is that agents involved in the making of one module do not observe the information relevant for the production of other modules. To improve efficiency, the firm can establish communication channels between agents who work on different modules. An IBM engineer who works on the disk drive does not directly observe the factors relevant to someone who works on the processor. But IBM can require the former to meet with the latter and learn about his local information. The problem is that such communication does not come for free. Even in the age of ever-evolving communication technologies, explaining the issues one faces, and understanding those faced by others, takes time and energy. Given this trade-off between the efficiency of decision-making and the cost of communication, the principal decides whom each agent should tell his state to. In terms of Figure 1, the principal takes as given

the production network in the left panel and designs the optimal communication network on the right by placing directed links between agents, taking into account that each comes at an exogenous cost but improves the agents' decision making.²

The challenge in designing an optimal communication network is the abundance of possibilities and absence of any apparent way to order them. To see how the firm can overcome this challenge, it is useful to start by asking when it should add a single, directed link to an existing communication network. The cost of establishing such a link is the time and energy it takes the sender to explain his information to the receiver, which we take to be exogenous. The benefit is that learning the sender's information allows the receiver to coordinate his decision more closely with the sender's, which, in turn, allows the sender to adapt his decision more closely to his state. Crucially, we find this benefit is independent of what the receiver, or any other agent, knows about any other state. Because of this independence, the problem of designing an optimal communication network can be broken into a number of independent subproblems. It is sufficient for the firm to consider each agent in turn and ask whom this agent should tell about his state. This separability result is what allows us to characterize optimal communication networks and establish our central result.

Our central result is that modular production gives rise to communication hierarchies. We fully characterize the optimal communication network and show that it consists of two nested hierarchies, one that determines whom each agent sends his information to and another that determines whom he receives information from. An agent with a higher sender rank sends his information to any agent that a lower-ranked agent sends his information to. And an agent with a higher receiver rank receives information from any agent that a lower-ranked agent receives information from.

An agent's rank in either hierarchy depends critically on the cohesiveness of his module, which is increasing in the number of decisions that are part of the module and the need for coordination between them, and which is decreasing in the "degree of coupling," the need for coordination across modules. The more cohesive an agent's module is, the more important it is that the agent learns about the local conditions in other modules and that agents working on other modules learn about his.

Receiver rank is fully determined by module cohesion. Agents working on the same module have the same receiver rank; they either all learn about a decision in another module or none of them do. The same is not true for sender rank, which can vary across agents working on the same module.

²Our focus on the trade-off between the efficiency of decision-making and the cost of communication is in line with Kenneth Arrow's discussion of the design of optimal communication within firms in Arrow (1974), where he observes that: "*Since information is costly, it is clearly optimal, in general, to reduce the internal transmission... That is, it pays to have some loss in value for the choice of terminal act in order to economize on internal communication channels. The optimal choice of internal communication structures is a vastly difficult question.*"

The reason is that the benefit of sharing information does not only depend on module cohesion but also on the variability of the sender’s local information and the need to adapt his decision to it. If an agent’s module is very cohesive, it is important he learns about other modules to enable them to adapt their decisions to their information. But if his own information is predictable, or the need to adapt his decision to his information is low, it may not be important for those working on other modules to learn his information. Sender and receiver ranks need not be perfectly correlated. Agents who hear the most may not be those who speak the most.

We apply our characterization to explore the notion that communication links should simply mirror technological interdependency, that “*we should expect to see a very close relationship...between a network graph of technical dependencies within a complex system and network graphs of organizational ties showing communication channels*” (Colfer and Baldwin 2016, p.713) This notion has a long history in management and related fields, where it is known as the *Mirroring Hypothesis* (Thompson (1967) and, for a discussion of the literature, Colfer and Baldwin (2016)). In our context, the Mirroring Hypothesis predicts communication within but not across modules. We show that for the Mirroring Hypothesis to hold, conditions have to be just right. There cannot be too many modules, and none of them can have too many members or require too much coordination among them. Otherwise the organization benefits from communication across modules. At the same time, none of the modules can have too few members and require too little coordination among them, which would negate the need for communication even within modules.

Our characterization result also speaks to the possibility that firms engage in “partial mirroring” by limiting across-module communication to clusters of modules. The result implies that if such clustering is optimal, there is at most one cluster of modules whose agents engage in across-module communication. Moreover, this cluster is made up of the most cohesive modules. Agents working on less cohesive modules are outside of the cluster and do not communicate across modules. In this case, the optimal communication network, therefore, has a core-periphery structure, which is a pervasive structure among social and communication networks (Herskovic and Ramos 2020). The model shows that this common structure can be an optimal way to organize a modular firm, but only if the core is made up of the modules with the largest number of decisions and/or the highest need for coordination among them.

2 Related Literature

To the best of our knowledge, there is no literature in economics on the economic implications of modular production. From a technical perspective, our paper belongs to the small but growing

literature on centralized network design. In an early paper in this literature, Baccara and Bar-Isaac (2008) explore the optimal design of a network among members of a criminal organization in which more links facilitate cooperation but also leave the organization more vulnerable to attack by law enforcement. The trade-off between the efficiency of interactions among members of a network and its increased vulnerability to attacks by outsiders is also at the center of Goyal and Vigier (2014), who are motivated by the optimal design and defense of computer networks. Both papers are quite far from ours, not only in terms of motivation but also modeling.

Much closer to us is Calvó-Armengol and de Martí (2009), who consider an organization in which each agent's payoff depends on how well his decision is adapted to a common state and coordinated with the other agents' decisions. A key feature of their model is that decisions enter the production function symmetrically: each has the same need for adaptation and is connected to all the others, with the need for coordination being the same between any two decisions. As such, their production network is complete, which precludes production from being modular. In our setting, in which the principal can add an arbitrary number of communication links at a given per link cost, such a production function gives rise to optimal communication networks that are bang-bang: each agent either tells his state to all the others or to none of them. In their setting, instead, the principal can add communication links at no cost up to an exogenously given cap. They show that, if the need for coordination is sufficiently small, and the degree of uncertainty sufficiently high, an optimal network then maximizes a span index that they define.

Herskovic and Ramos (2020) also consider a setting in which each agent's payoff depends on how well his decision is adapted to a single state and coordinated with the other decisions, and in which all decisions enter the production function symmetrically. The key difference between their model and both Calvó-Armengol and de Martí (2009) and ours is that the communication network is not designed by a principal but formed by the agents' decentralized decisions of whom to communicate with. Their paper, therefore, belongs to the large literature on endogenous network formation that started with Jackson and Wolinsky (1996) and Bala and Goyal (2000), rather than the literature on centralized network design that ours belongs to. They show that, in spite of agents' decisions being identical *ex ante*, the network they form is hierarchical, with agents in a given tier having their signals being observed by those in the lower tiers.

In the endogenous network formation literature, the paper whose setting is closest to ours is Calvó-Armengol, de Martí, and Prat (2015). They, too, consider an organization whose agents face a trade-off between adaptation and coordination. Like us, though, they assume that each agent is adapting his decision to an independent state and, crucially, allow for decisions to differ in their

needs for adaptation and coordination. Even though they allow for differing coordination needs, however, they do not assume that it has a non-overlapping community structure, and thus do not explore modular production. Moreover, they assume that each agent decides independently how much effort to put into communicating with each of the other agents, which is what places them in the endogenous network formation literature. In contrast, we allow the principal to decide who each agent communicates with. Their main result provides a characterization of how, in such a setting, an agent's decision is influenced by the signals received by others.

A shared feature between all the above papers on adaptation and coordination, and ours, is that the agents' payoff functions are quadratic, and their actions are continuous and exhibit strategic complementarities. As such, they all build on the literature on quadratic games on networks that started with Ballester, Calvó-Armengol, and Zenou (2006). In recent contributions to this literature, Bergemann, Heumann, and Morris (2017) and Golub and Morris (2017) characterize optimal decision-making for general information and network structures. We draw on their results to determine the agents' decision-making for given communication networks. Our focus, though, is not on the agents' decision-making but on the prior stage in which the principal designs the communication network, taking as given that agents will make their decisions optimally.

Apart from the literature on networks, our motivation and application places us firmly in the literature on organizational economics and, in particular, team theory. Starting in the 1950s, team theory explores the optimal design of organizations when agents share the same goal, but cognitive constraints make communication costly (for an early treatment see Marschak and Radner (1972) and for recent surveys see Garicano and Prat (2011) and Garicano and van Zandt (2013)). In this literature, a closely related paper is Dessein and Santos (2006), who were the first to explore how the trade-off between adaptation and coordination shapes the internal organization of firms. In their setting, decisions enter the production function symmetrically, and the principal does not design a communication network. Instead, they allow for each agent to make multiple decisions and assume the same quality of communication between any pair of agents. They show that, in such a setting, more uncertainty about the environment increases the optimal number of decisions per agent, while the effect of an improvement in the quality of communication on specialization is non-monotonic.

We also relate to Dessein, Galeotti, and Santos (2016), who build on Dessein and Santos (2006) by endogeneizing communication while taking the allocation of decisions as given. Decisions differ in their needs for adaptation but the need for coordination is the same for any pair of decisions. As such, the production network is complete, and production is not modular. As in our model,

the principal designs the firm's communication structure before agents learn their information and make their decisions. In contrast to our model, each agent's communication is public, leaving all the other agents equally well, or poorly, informed about him. The paper shows that if the total amount of time that agents have to learn about others is limited, the principal finds it optimal to have them spend all their time learning about a small number of core agents, while staying largely ignorant about the others.

Even though, to our knowledge, there is no literature on modularity in economics, there is a large literature on this topic in management and related fields, as well as in computer science. As noted earlier, Simon (1962) observed that complex systems are often made up of modules and argued that this modular design facilitates adaptation. A similar point was made by Alexander (1964), who argued that a modular system design accelerates adaptation by allowing the system to adapt module by module. In computer science, Parnas (1972) argued that a modular software design allows for faster programming by enabling different teams to work on different program modules in parallel, and explored criteria to best decompose a program into modules.

Our paper connects to a related literature that takes the modular design of products as given and explores its implications for the organization of production. A common argument in this literature is the Mirroring Hypothesis we mentioned in the introduction, which posits that the organization of a firm, and specifically its internal communication structure, ought to mirror the modular nature of its production function. A firm that makes a modular product, in other words, should see intense communication within modules but not across (see, in particular, Thompson (1967), Henderson and Clark (1990), Sanchez and Mahony (1996) and, for a survey, see Baldwin and Colfer (2016)). Langlois and Robertson (1992) observed that modular production might not only affect the internal organization of firms but also their boundaries and, through this channel, the structure of industries. Baldwin and Clark (2000) document these dynamics in the context of IBM and the computer industry, and provide an exhaustive discussion of modular production and its organization.

A related literature reverses the causality of the Mirroring Hypothesis and argues that the design of products reflects the organization of the firms that developed them. In this view, a modular organization has a tendency to develop modular products. In computer science, this view is known as *Conway's Law*, named after Melvin Conway who observed that "*To the extent that an organization is not completely flexible in its communication structure, that organization will stamp out an image of itself in every design it produces*" Conway (1968, p.30).

3 Model

A firm consists of one principal and N agents. All parties are risk neutral and care only about the firm's profits. There are no incentive conflicts.

Production. Each agent $i \in \mathcal{N}$ makes a decision $d_i \in [-D, D]$ that is associated with a state $\theta_i \in [-D, D]$, where $\mathcal{N} = \{1, \dots, N\}$ is the set of agents and D is a large but finite scalar. Output depends on how well each decision is adapted to its associated state and coordinated with the other decisions. Specifically, output is given by

$$r(d_1, \dots, d_n) = \sum_{i=1}^N \left[-d_i^2 + 2a_i d_i \theta_i + \sum_{j=1}^N p_{ij} d_i d_j \right], \quad (1)$$

where $a_i > 0$ captures the importance of adapting decision d_i to its state θ_i , $p_{ij} \geq 0$ represents the need to coordinate decisions $d_i \neq d_j$, and where $p_{ii} = 0$. The need for coordination is symmetric, that is, $p_{ij} = p_{ji}$. The interactions between decisions can, therefore, be represented by an undirected network, which we summarize in an $N \times N$ matrix \mathbf{P} with entries p_{ij} . We assume that $\sum_{j=1}^N p_{ij} < 1$ for all $i \in \mathcal{N}$, which ensures that equilibrium decisions exist. Finally, we normalize the price of the product to one so that output (1) also represents revenue.

Modules. Each decision, and its associated state and agent, belongs to a “module” \mathcal{M}_m for $m \in \{1, \dots, M\}$, which is a set of $n_m \geq 1$ such decisions. The function $m(i)$ gives the module $\mathcal{M}_{m(i)}$ that decision d_i belongs to. For expositional convenience we adopt the convention that the first decision d_1 , and its associated state and agent, belong to module \mathcal{M}_1 .

The need for coordination is stronger between two decisions within the same module than between two decisions in different modules. Specifically, the need for coordination between any two decisions d_i and $d_j \neq d_i$ is given by $p_{ij} = t \geq 0$ if they belong to different modules and, abusing notation slightly, it is given by $p_{ij} = p_m \geq t$ if they belong to the same module \mathcal{M}_m . The parameter t , therefore, captures the degree of coupling between modules, while the parameter p_m captures the need for coordination within module \mathcal{M}_m .

Information. Each state θ_i is independently drawn from a distribution with zero mean and variance σ_i^2 , for any $i \in \mathcal{N}$. The realization of state θ_i is privately observed by agent i and the other agents in his module $\mathcal{M}_{m(i)}$. All other information is public.

Before the states are realized, the principal can place directed communication links between any two agents. Each such link comes at a cost $\gamma > 0$, which captures the resources involved in communication. If the principal places a communication link from agent i to agent j , agent i tells

j the realization of his state θ_i . Communication, therefore, takes the form of a directed network, which we summarize in the $N \times N$ matrix \mathbf{C} . Entry c_{ij} is equal to one if agent i tells agent j about his state, or agent j observes θ_i directly, and it is equal to zero otherwise. Row \mathbf{C}_i then summarizes the agents who learn θ_i and column $\mathbf{C}_{(j)}$ summarizes the states agent j learns about.

Organization. The principal's problem is to design the optimal communication network that maximizes expected revenue net of communication costs, that is, to solve

$$\max_{\mathbf{C}} \mathbb{E}[r(d_1, \dots, d_N) | \mathbf{C}] - \gamma \sum_{i=1}^N (\mathbf{C}_i \mathbf{1} - n_{m(i)}), \quad (2)$$

subject to $c_{ij} = 1$ for all $i, j \in \mathcal{N}$ such that $\mathcal{M}_{m(i)} = \mathcal{M}_{m(j)}$, where $\mathbf{1}$ is an $N \times 1$ vector of ones.

Timing. After the principal designs the communication network, agents learn their states and tell them to other agents as specified in the network. Next, the agents simultaneously make their decisions, payoffs are realized, and the game ends. The solution concept we use is Perfect Bayesian Equilibrium.

We discuss the various key assumptions, such as the assumptions that agents do not re-transmit information they receive and that their decisions are not distorted by incentive conflicts, in Section 8, after solving the model in the next three sections and applying it to the Mirroring Hypothesis in Section 7.

4 Decision-Making

We solve the game by first determining equilibrium decisions for any given communication network and then characterizing optimal communication networks. After agents have observed their states and communicated with each other, they make the decisions that solve

$$\max_{d_i} \mathbb{E}[r(d_1, \dots, d_N) | \mathbf{C}_{(i)}] \quad \text{for all } i \in \mathcal{N},$$

where $r(d_1, \dots, d_N)$ is revenue (1) and where $\mathbf{C}_{(i)}$ is the i 's column of the communication matrix \mathbf{C} that summarizes the states agent i is informed about. The best response functions that follow from these optimization problems are given by

$$d_i = a_i \theta_i + \sum_{j=1}^N p_{ij} \mathbb{E}[d_j | \mathbf{C}_{(i)}]. \quad (3)$$

Each agent's best-response, therefore, is the weighted sum of his state and the decisions he expects the other agents to make, where the weight on his own state is a_i and the weight on the decision he

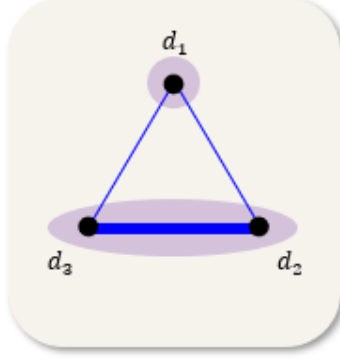


Figure 3: Production network P^e . The width of the edges indicates the need for coordination and the blue shaded areas highlight the modules.

expects agent j to make is p_{ij} . To solve the system of best responses, note that $(\text{diag } \mathbf{C}_j) \mathbf{P} (\text{diag } \mathbf{C}_j)$ is the subgraph of the production network that consists of the nodes whose agents know state θ_j , as well as all the links between them. We can then state the following lemma.

LEMMA 1. *Equilibrium decisions are unique and given by*

$$d_i^* = \sum_{j=1}^N a_j \omega_{ij}(\mathbf{C}_j) \theta_j \text{ for all } i \in \mathcal{N}, \quad (4)$$

where $\omega_{ij}(\mathbf{C}_j)$ denotes the ij th entry of $(\mathbf{I} - (\text{diag } \mathbf{C}_j) \mathbf{P} (\text{diag } \mathbf{C}_j))^{-1}$.

The lemma shows that agent i 's equilibrium decision d_i^* is the weighted sum of all states, where the weight on state θ_j is given by a_j , the importance of adapting decision d_j to θ_j , times $\omega_{ij}(\mathbf{C}_j)$, the ij th entry of $(\mathbf{I} - (\text{diag } \mathbf{C}_j) \mathbf{P} (\text{diag } \mathbf{C}_j))^{-1}$. This latter object has a natural interpretation in terms of walks on the production network, but before providing it, we pause briefly to review the notion of walks and their values.

A “walk” between d_i and d_j on the production network is a sequence of links that lead from d_i to d_j . Each link between two decisions in this sequence is associated with a discount factor, which is given by the need for coordination between them. The “value” of a walk is the product of these discount factors. As an example, consider the production network P^e in Figure 3, where the superscript stands for “example.” In this case, d_1 to d_2 to d_3 constitutes a walk from d_1 to d_3 whose value is given by $p_{12}p_{23}$. Standard arguments imply that the ij th entry of $(\mathbf{I} - \mathbf{P}^e)^{-1}$ is the sum of the values of all walks from d_i to d_j on the production network P^e .

In light of this discussion, the ij th entry of $(\mathbf{I} - (\text{diag } \mathbf{C}_j) \mathbf{P} (\text{diag } \mathbf{C}_j))^{-1}$ in agent i 's equilibrium decision rule (4) represents the value of all walks from node i to node j on the subgraph of the

production network that consists only of decisions made by agents who know state θ_j . If agent i does not know θ_s , for instance, d_i is not part of this subgraph, and so $\omega_{is}(\mathbf{C}_s) = 0$. Agent i puts no weight on θ_s , as one would expect. If, instead, θ_s is public, the subgraph encompasses the entire production network and the weight agent i puts on θ_s is the value of all walks from d_i to d_s on the production network \mathbf{P} . Note that this is the case no matter what the agents know about the other states. This result reflects a general implication of the lemma that will be important for what follows: for a given production network, the weight agent i puts on state θ_s depends only on who knows θ_s and not on what agent i , or any other agent, knows about any other state.

To get an intuition for the equilibrium decision rule in Lemma 1, recall the production network \mathbf{P}^e in Figure 3 and consider how much weight agent 1 puts on his own state. If only agent 1 observes his state, agents 2 and 3 cannot put any weight on θ_1 . As a result, agent 1 faces a trade-off between adapting his decision to θ_1 and coordinating it with d_2 and d_3 . Lemma 1 shows that if agent 1 resolves this trade-off optimally, the weight he puts on his own state is given by

$$a_1 \omega_{11}^e((1, 0, 0)) = a_1,$$

where the superscript again stands for “example,” and where the equality follows from the fact that the value of all walks from d_1 to d_1 that go only through d_1 is one.

If agent 2 learns θ_1 , his decision will put some weight on this state, which relaxes agent 1’s trade-off between adapting his decision to θ_1 and coordinating it with d_2 . As a result, agent 1 increases the weight he puts on his state to

$$a_1 \omega_{11}^e((1, 1, 0)) = a_1 \left(1 + \frac{t^2}{1 - t^2} \right),$$

where the fraction on the right-hand side is the value of all walks from d_1 back to d_1 that go through d_1 and d_2 .

Finally, if agent 3 also learns θ_1 , the weight agent 1 puts on his state increases further to

$$a_1 \omega_{11}^e((1, 1, 1)) = a_1 \left(1 + \frac{t^2}{1 - t^2} + \frac{t^2}{1 - t^2} + \frac{2t^2(t^2 + p_{23})}{(1 - t^2)(1 - p_{23} - 2t^2)} \right),$$

where the second fraction is the value of all walks from d_1 back to d_1 that only go through d_1 and d_3 , and the third fraction is the value of all walks from d_1 back to d_1 that go through all three decisions. The intuition for the second fraction is the same as for the first: if agent 3 learns θ_1 , his decision will put some weight on this state. This relaxes the trade-off agent 1 faces between adapting d_1 to θ_1 and coordinating it with d_3 , inducing him to put more weight on his state. The additional effect that is captured by the third fraction is that by making his decision vary with θ_1 ,

agent 3 also relaxes the trade-off agent 2 faces between coordinating d_2 with both d_1 and d_3 . As a result, agent 2 increases the weight his decision puts on θ_1 , which allows agent 1 to do the same.

The example illustrates three general features of the agents' decision-making that matter for what follows. First, if only agent i observes θ_i , he puts weight a_i on his state. The parameter a_i , therefore, captures the degree of autonomous adaptation. Second, allowing additional agents to observe θ_i scales the weight agent i puts on his state to $\omega_{ii}(\mathbf{C}_i) a_i \geq a_i$. It does so since allowing additional agents to learn θ_i enables them to coordinate their decisions with each other and with d_i , which, in turn, allows agent i to adapt his decision to θ_i without sacrificing as much coordination. The total weight that agent i puts on his state is, therefore, the product of the "coordination multiplier" $\omega_{ii}(\mathbf{C}_i)$ and the degree of autonomous adaptation a_i . Finally, the increase in the weight agent i puts on his state if agent $j \neq i$ learns θ_i is larger, if agent $k \neq i, j$ knows θ_i . We summarize these properties in the following corollary.

COROLLARY 1. *The weight $a_i \omega_{ii}(\mathbf{C}_i)$ that agent i 's decision d_i^* puts on his state θ_i satisfies $\omega_{ii}(\mathbf{I}_i) a_i = a_i$ and is increasing and supermodular in \mathbf{C}_i .*

Having characterized equilibrium decision-making by the agents, we next turn to the principal's problem.

5 The Principal's Problem

The principal's problem is to design the communication network that maximizes expected profits, taking into account that agents make decisions according to (4). It is useful to start by rewriting revenue (1) as

$$r(d_1, \dots, d_N) = \sum_{i=1}^N a_i d_i \theta_i - \sum_{i=1}^N d_i \left(d_i - a_i \theta_i - \sum_{j=1}^N p_{ij} d_j \right).$$

Substituting in the equilibrium decision rules (4), this simplifies to

$$r(d_1^*, \dots, d_N^*) = \sum_{i=1}^N a_i d_i^* \theta_i + \sum_{i=1}^N \sum_{j=1}^N p_{ij} d_i^* (d_j^* - \mathbb{E}[d_j^* | \mathbf{C}_{(i)}]). \quad (5)$$

In the proof of the next lemma we show that the second term on the right-hand side is zero in expectation, which implies the following result.

LEMMA 2. *Under equilibrium decision-making, expected revenue is given by*

$$R(\mathbf{C}) \equiv \mathbb{E}[r(d_1^*, \dots, d_N^*)] = \sum_{i=1}^N a_i \text{Cov}(d_i^*, \theta_i), \quad (6)$$

where $\text{Cov}(d_i^*, \theta_i) = a_i \sigma_i^2 \omega_{ii}(\mathbf{C}_i)$.

The lemma shows that expected revenue boils down to how well each decision is adapted to its associated state. For expositional convenience, we interpret $a_i \text{Cov}(d_i^*, \theta_i)$ as the expected revenue generated by agent $i \in \mathcal{N}$ and denote it by

$$R_i(\mathbf{C}_i) \equiv a_i \text{Cov}(d_i^*, \theta_i) = \sigma_i^2 a_i^2 \omega_{ii}(\mathbf{C}_i).$$

The key property of agent i 's expected revenue is that it depends on \mathbf{C}_i but not \mathbf{C}_{-i} , on who knows θ_i but not on what agent i , or any other agent, knows about any other state. An additional agent learning θ_i increases agent i 's coordination multiplier $\omega_{ii}(\mathbf{C}_i)$, and thus the weight $a_i \omega_{ii}(\mathbf{C}_i)$ he puts on his state, as well as the expected revenue $\sigma_i^2 a_i^2 \omega_{ii}(\mathbf{C}_i)$ he generates. In contrast, agent i , or any other agent, learning any other state does not affect $\omega_{ii}(\mathbf{C}_i)$, and thus leaves the weight agent i puts on his state, and the revenue he is expected to generate, unchanged.

This property of expected revenue is key because it implies that the principal's problem is separable. Instead of solving the overall problem (2) head on, the principal can consider each agent in isolation and ask whom this agent should tell about his state. The answer to whom agent $i \in \mathcal{N}$ should tell about θ_i is independent of whom any other agent should tell about his own state. We, therefore, have our first main result.

PROPOSITION 1. *An optimal communication network solves the principal's problem (2) if and only if it solves the N independent subproblems*

$$\max_{\mathbf{C}_i} R_i(\mathbf{C}_i) - \gamma (\mathbf{C}_i \mathbf{1} - n_{m(i)}) \text{ for all } i \in \mathcal{N}. \quad (7)$$

This separability result greatly facilitates the principal's quest for optimal communication networks. We can further simplify the problem by recalling that agent i 's coordination multiplier $\omega_{ii}(\mathbf{C}_i)$ is supermodular. This property implies that, whenever it is optimal for agent i to tell agent j about his state, it must also be optimal for him to tell the other agents in agent j 's module $\mathcal{M}_{m(j)}$. The principal's problem, therefore, reduces to which modules each agent should tell about his state.

To state the principal's problem in these terms, let \mathbf{G} denote the matrix that specifies the modules whose agents are told about each state, and let \mathbf{F} denote the conversion matrix that specifies the module each agent belongs to. Specifically, \mathbf{G} is an $N \times M$ matrix in which entry g_{im} is equal to one if agent i tells agents in module \mathcal{M}_m about his state θ_i , and g_{im} is equal to zero if he does not. Moreover, $g_{im(i)} = 1$, which reflects that agents observe the states in their own module without having to be told about them. The conversion matrix \mathbf{F} , in turn, is an $M \times N$ matrix in

which entry f_{mi} is equal to one if agent i belongs to module \mathcal{M}_m and zero otherwise. We can then restate the principal's problem (7) as

$$\max_{\mathbf{G}_i} R_i(\mathbf{G}_i \mathbf{F}) - \gamma \left(\sum_{m=1}^M g_{im} n_m - n_{m(i)} \right) \text{ for all } i \in \mathcal{N}. \quad (8)$$

Finally, supermodularity of $\omega_{ii}(\cdot)$, together with the linearity of the communication costs, imply that the subproblems (8) are also supermodular. For any given parameter values, the principal's problem can, therefore, be solved using standard algorithms that maximize supermodular functions in polynomial time (see, for instance, chapter 10.2 in Murota (2003)). Our goal, though, is to solve the problem analytically, and we do so in the next section.

6 Optimal Communication Networks

The separability result in Proposition 1 allows us to solve the principal's problem of designing optimal communication networks by considering each agent in isolation and asking whom he should tell about his state. To economize on notation, we focus on agent 1. Once we know who agent 1 should tell about his state, we can apply the answer to all the other agents, and thus solve the principal's overall problem (8).

To this end, consider the expected revenue $R_1(\cdot)$ agent 1 generates if his state is known to the agents in his own module \mathcal{M}_1 and to those in some arbitrary set of other modules. Since the naming of modules is immaterial, there is no loss in denoting the modules whose agents know θ_1 by $\mathcal{M}_1, \dots, \mathcal{M}_\ell$ for $\ell \in \{1, \dots, M\}$. We can then define $\mathbf{G}_1(\ell)$ as the $1 \times M$ row vector that specifies the modules whose agents know θ_1 and $\mathbf{C}_1(\ell) = \mathbf{G}_1(\ell) \mathbf{F}$ as the corresponding $1 \times N$ row vector of the communication matrix that specifies the agents who know θ_1 . The next lemma uses this notation to express agent 1's expected revenue.

LEMMA 3. *Suppose agent 1's state θ_1 is known to all the agents in modules $\mathcal{M}_1, \dots, \mathcal{M}_\ell$, for $\ell \in \{1, \dots, M\}$, and none of the agents in other modules. Agent 1's expected revenue is then given by*

$$R_1(\mathbf{C}_1(\ell)) = a_1^2 \sigma_1^2 \left(\frac{1 - (n_1 - 2)p_1}{(1 + p_1)(1 - (n_1 - 1)p_1)} + \frac{t^2 x_1^2 \sum_{m=2}^{\ell} n_m x_m}{(1 - t n_1 x_1) \left(1 - t \sum_{m=1}^{\ell} n_m x_m\right)} \right), \quad (9)$$

where

$$x_m \equiv \frac{1}{1 - (n_m - 1)p_m + n_m t} \text{ for } m = 1, \dots, M.$$

The object x_m in the lemma captures the “cohesion” of module \mathcal{M}_m , which is increasing in its size and the need for coordination among its members, and decreasing in the degree of coupling t . A module is, therefore, more cohesive if coordination among its members is relatively more important than coordination between its members and those in other modules.³ Both module cohesion x_m and its scaled analog $n_m x_m$ play a key role in who agent 1 should tell about his state.

In line with our discussion in the previous section, Lemma 3 shows that agent 1’s expected revenue is the product of $a_1^2 \sigma_1^2$ and the coordination multiplier $\omega_{11}(\mathbf{C}_1(\ell))$, which itself is the sum of two terms. The first term is the value of all walks from node 1 back to itself that only go through nodes in \mathcal{M}_1 . Because these walks only go through module \mathcal{M}_1 , their value depends only on its characteristics p_1 and n_1 . The second term, in turn, is the value of the additional walks that also go through the other modules $\mathcal{M}_2, \dots, \mathcal{M}_\ell$. Notice that the value of these additional walks depends on the characteristics of \mathcal{M}_1 only through x_1 and $n_1 x_1$, and that it depends on the characteristics of the other informed modules only through the sum of their $n_m x_m$ terms. The revenue generated by agent 1, for instance, is the same whether he tells his state to agents in one module with $n_2 x_2 = 10$ or to agents in ten modules with $n_2 x_2 = \dots = n_{11} x_{11} = 1$. This property of expected revenue is important for what follows.

Having derived agent 1’s expected revenue, we can now determine when the principal benefits from having him tell his state, not just to agents in modules $\mathcal{M}_2, \dots, \mathcal{M}_\ell$, but also to those in either only $\mathcal{M}_{\ell+1}$ or in both $\mathcal{M}_{\ell+1}$ and $\mathcal{M}_{\ell+2}$. We will see below that doing so delivers a characterization of optimal communication.

To this end, suppose that, initially, there are at least two modules whose agents do not know θ_1 , that is, $\ell \leq M - 2$, and that the principal expands the set of informed agents by having agent 1 also tell his state to the agents in module $\mathcal{M}_{\ell+1}$. She benefits from doing so as long as the per node marginal revenue it generates is larger than the marginal cost, that is, as long as

$$\frac{1}{n_{\ell+1}} (R_1(\mathbf{C}_1(\ell+1)) - R_1(\mathbf{C}_1(\ell))) \geq \gamma.$$

Using (9), the per node marginal revenue is given by

$$a_1^2 \sigma_1^2 \frac{t^2 x_1^2 x_{\ell+1}}{\left(1 - t \sum_{m=1}^{\ell} n_m x_m\right) \left(1 - t \sum_{m=1}^{\ell+1} n_m x_m\right)}, \quad (10)$$

where the fraction is the change in the coordination multiplier $\omega_{11}(\mathbf{C}_1(\ell))$. Notice that the change in the coordination multiplier is increasing in x_1 and $n_1 x_1$ and thus in n_1 and p_1 . The more

³There are different notions and formal definitions of cohesion in the sociology and economics literatures. Our definition is close to that in Morris (2000). Applied to our setting, his definition of the cohesion of module \mathcal{M}_m is $(n_m - 1) p_m / [(n_m - 1) p_m + (N - n_m) t]$.

members module \mathcal{M}_1 has, or the more important it is to coordinate their decisions, the more the principal gains from agents in an additional module learning θ_1 . Beyond the characteristics of agent 1's module, the change in the coordination multiplier is increasing in the number of informed modules, which reflects the supermodularity of $\omega_{11}(\cdot)$. And it is increasing in p_m and n_m for all $m = 2, \dots, \ell + 1$, which reflects the fact that the value of the additional walks that are created by telling agents in another module about a state is increasing in the size of the module and the need for coordination among them.

Notice also that while the change in the coordination multiplier depends only on characteristics of the different modules, per node marginal revenue also depends on $a_1^2 \sigma_1^2$, which is specific to decision d_1 and may be different from the corresponding values for the other decisions in the same module. The term $a_1^2 \sigma_1^2$ captures the revenue agent 1 would generate if he were the only one who knew θ_1 and had to adapt his decision to his state without others coordinating their decisions with his. As such, $a_i^2 \sigma_i^2$ is the ‘‘value of autonomous adaptation’’ that captures the importance of adapting decision $i \in \mathcal{N}$ to its state. The higher the value of autonomous adaptation $a_1^2 \sigma_1^2$ is, the higher is the per node marginal revenue (10) for any given module characteristics.

Next it is useful to compare the per node marginal revenue that is generated when agent 1 tells his state to agents in one more module with that when he tells it to those in two modules. Suppose the principal extends the communication network by having agent 1 tell his state to agents in both $\mathcal{M}_{\ell+1}$ and $\mathcal{M}_{\ell+2}$. Again using (9), the per node marginal revenue of doing so

$$\frac{1}{n_{\ell+1} + n_{\ell+2}} (R_1(\mathbf{C}_1(\ell + 2)) - R_1(\mathbf{C}_1(\ell)))$$

is given by

$$a_1^2 \sigma_1^2 \frac{1}{n_{\ell+1} + n_{\ell+2}} \frac{t^2 x_1^2 (n_{\ell+1} x_{\ell+1} + n_{\ell+2} x_{\ell+2})}{\left(1 - t \sum_{m=1}^{\ell} n_m x_m\right) \left(1 - t \sum_{m=1}^{\ell+2} n_m x_m\right)}.$$

Subtracting (10) from this expression, we have that the difference in the per node marginal revenues

$$\frac{1}{n_{\ell+1} + n_{\ell+2}} (R_1(\mathbf{C}_1(\ell + 2)) - R_1(\mathbf{C}_1(\ell))) - \frac{1}{n_{\ell+1}} (R_1(\mathbf{C}_1(\ell + 1)) - R_1(\mathbf{C}_1(\ell)))$$

is equal to

$$a_1^2 \sigma_1^2 \frac{n_{\ell+2} t^2 x_1^2}{n_{\ell+1} + n_{\ell+2}} \frac{(x_{\ell+2} - x_{\ell+1}) \left(1 - t \sum_{m=1}^{\ell+1} n_m x_m\right) + t(n_{\ell+1} + n_{\ell+2}) x_{\ell+1} x_{\ell+2}}{\left(1 - t \sum_{m=1}^{\ell} n_m x_m\right) \left(1 - t \sum_{m=1}^{\ell+1} n_m x_m\right) \left(1 - t \sum_{m=1}^{\ell+2} n_m x_m\right)}.$$

The key property of this expression is that it is positive if $x_{\ell+2} \geq x_{\ell+1}$. If module $\mathcal{M}_{\ell+2}$ is more cohesive than module $\mathcal{M}_{\ell+1}$, the per node marginal revenue of telling agents in both $\mathcal{M}_{\ell+1}$ and

$\mathcal{M}_{\ell+2}$ about the state is larger than that of telling only those in $\mathcal{M}_{\ell+1}$. It can then never be optimal for agent 1 to tell his state to agents in $\mathcal{M}_{\ell+1}$ but not to those in $\mathcal{M}_{\ell+2}$. It may be optimal for agent 1 to tell agents in both modules, neither, or only those in $\mathcal{M}_{\ell+2}$, but telling only the agents in $\mathcal{M}_{\ell+1}$ cannot be optimal. This property, in turn, implies our main result: optimal communication follows a threshold rule in which agent 1 tells his state to agents in any module \mathcal{M}_m whose cohesion x_m is above a threshold, and to none in those whose value is below. Since there is nothing special about agent 1, the same applies to any other agent. We can then characterize the solution to the principal's problem.

PROPOSITION 2. *Optimal communication is characterized by N thresholds $\lambda_i \geq 0$, one for each agent $i \in \mathcal{N}$. Given any two agents i and j who belong to different modules, agent i tells agent j about his state if and only if the cohesion of agent j 's module is above agent i 's threshold, that is, if and only if $x_{m(j)} \geq \lambda_i$. The threshold λ_i is increasing in marginal communication costs γ and decreasing in the value of autonomous adaptation $a_i^2 \sigma_i^2$, the need to coordinate the decisions within agent i 's module $p_{m(i)}$, and the size of his module $n_{m(i)}$.*

To illustrate the proposition for agent 1, it is convenient to label modules $\mathcal{M}_2, \dots, \mathcal{M}_\ell$ in decreasing order of their values of x_m , so that \mathcal{M}_2 denotes the module, other than \mathcal{M}_1 , with the highest value of x_m , and \mathcal{M}_M denotes the one with the smallest. In Figure 4, the blue curve is the piecewise linear extension of expected revenue $R_1(\mathbf{C}_1(\ell))$, which we denote by $\bar{R}_1(\mathbf{C}_1(\ell))$, and the red line is a continuous representation of communication costs $\sum_{m=1}^{\ell} n_m \gamma$. The changing curvature of expected revenue $\bar{R}_1(\mathbf{C}_1(\ell))$ reflects the countervailing economic forces at work. The supermodularity at the heart of the model pushes towards convexity while the modular structure of the production function pushes towards concavity. A reduction in γ favors telling agents in more modules about the state because it flattens the cost curve. And an increase in the value of autonomous adaptation $a_1^2 \sigma_1^2$ increases the per node marginal benefit (10), and thus steepens the expected revenue curve, as does an increase in the size of the agent 1's module n_1 or the need for coordination among its members p_1 .

The proposition implies that optimal communication gives rise to sender and receiver hierarchies. To see this clearly, focus again on agent 1 and consider the agents in modules \mathcal{M}_2 and \mathcal{M}_3 . The proposition shows that if \mathcal{M}_2 is more cohesive than \mathcal{M}_3 , agents in module \mathcal{M}_3 will only ever be told about θ_1 if those in module \mathcal{M}_2 also are. Moreover, since x_2 and x_3 do not depend on the characteristics of the sender's module \mathcal{M}_1 , agents in module \mathcal{M}_3 will only ever be told about *any* state that those in module \mathcal{M}_2 also are. Optimal communication, therefore, gives rise to a *receiver hierarchy*, in which a higher-ranked agent is told about all the states that a lower-ranked agent is

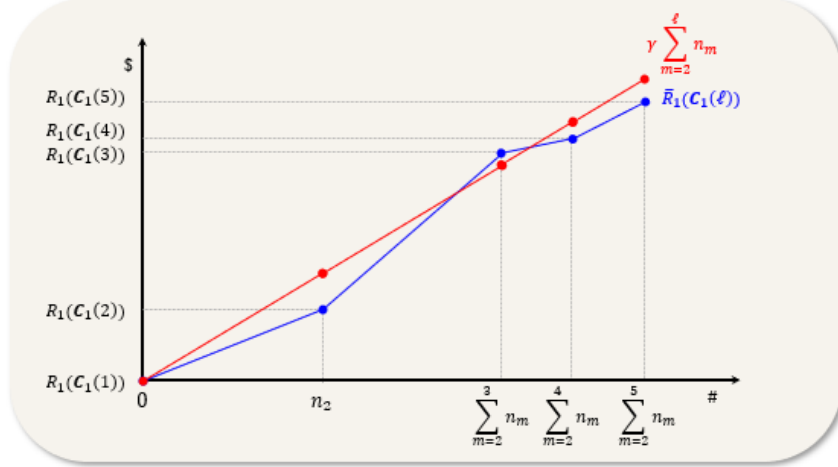


Figure 4: Determining the optimal communication network for agent 1 (drawn for parameter values $t = 0.01$, $n_1 = n_2 = n_3 = 5$, $n_4 = n_5 = 2$, $p_1 = p_2 = p_3 = 0.2$, $p_4 = p_5 = 0.1$, $a_1\sigma_1 = 1$).

told about, and possibly more. Agent i 's position in the ranking is fully determined by the cohesion $x_{m(i)}$ of the module he belongs to. If $x_{m(i)} \geq x_{m(j)}$, agent i outranks agent j .

COROLLARY 2. *Optimal communication gives rise to a receiver hierarchy among agents. For any agents $i, j, k \in \mathcal{N}$ who belong to different modules, if agent i 's module is more cohesive than agent j 's, then agent j is told about agent k 's state only if agent i also is.*

Notice that this result is about communication and not information per se. A higher-ranked agent is told about all the states that a lower-ranked agent is told about. But a lower-ranked agent may still have some information that a higher-ranked agent does not have. In particular, a lower-ranked agent observes his own state, and those in his module, directly and it may well be optimal for a higher-ranked agent to remain ignorant about those states.

The fact that cohesion x_m depends only on the characteristics of module \mathcal{M}_m , and not on those of any other modules, also implies a sender hierarchy in which a higher-ranked agent tells his state to all the agents that a lower-ranked agent does, and possibly others. The rank of agent i , though, does not depend on $x_{m(i)}$ but on λ_i , and it does so inversely. Agent i outranks agent j if $\lambda_i \leq \lambda_j$.

COROLLARY 3. *Optimal communication gives rise to a sender hierarchy among agents. For any agents $i, j, k \in \mathcal{N}$ who belong to different modules, if agent i 's threshold λ_i is smaller than agent j 's threshold λ_j , then agent j tells agent k about his state only if agent i also does.*

Even though an agent's rank in either hierarchy is increasing in the cohesion of his module, the positions need not coincide. Agent i may outrank agent j in one hierarchy but be outranked by

him in the other. The reason is that while agent i 's rank in the receiver hierarchy depends only the cohesion of the modules, his rank in the sender hierarchy also depends on the specific characteristics of his decision, as captured by the value of autonomous adaptation $a_i^2\sigma_i^2$. Suppose agent i 's module is very cohesive so that $x_{m(i)}$ is larger than the x_m of any other module. Agent i , and the other agents in his module, then reside on top of the receiver hierarchy, receiving any communication shared with any agents in any other module. They reside on top of the receiver hierarchy because their ignorance about other modules would hold those modules back from adapting their decisions more than the ignorance of agents in any other module would. At the same time, if $a_i^2\sigma_i^2$ is sufficiently small, λ_i is also smaller than the λ_j of any other agent $j \in \mathcal{N} \setminus i$, placing agent i at the bottom of the sender hierarchy. Even though his module is very cohesive, his ability to adapt his decisions to his state is just not very important. The agents who hear the most, therefore, might also speak the least.

A distinct feature of both hierarchies is that they are nested. This feature is in contrast to the properties of the knowledge hierarchies in Garicano (2000), and the literature that builds on his work, in which an agent tells his immediate boss about a problem he cannot solve, but does not tell the boss's superiors. The choice of whether to tell the boss's superior about the problem is left with the boss, who sometimes decides to do so and sometimes does not. Such narrow communication chains are not optimal in our setting. If an agent tells a superior about his state, he tells the superior's superiors. And if he hears from a subordinate, he hears from the subordinate's subordinates.

7 Application

Our result that the optimal organization of modular production is hierarchical contrasts with the Mirroring Hypothesis. As we discussed earlier, the Mirroring Hypothesis conjectures that the optimal way to organize modular production is to simply mirror the production function, to ensure intense communication within modules and accept sparse communication across.

The Boeing Company's experience with the 787 Dreamliner illustrates the Mirroring Hypothesis and why it may not always hold.⁴ The Dreamliner was designed to be modular precisely because it allowed Boeing to outsource the development and production of most modules to independent suppliers, many of which were scattered around the globe (see Figure 5). Suppliers delivered the finished modules to Boeing's factory in Everett, where its workers put them together with the

⁴This account is based on Peterson (2011) and Brown and Garthwaite (2016). See also Tadelis and Williamson (2012).

tail fin, the only major module still made by Boeing itself. To the extent that firm- and country boundaries hamper communication, this way of organizing the production of the Dreamliner is broadly in line with the Mirroring Hypothesis.

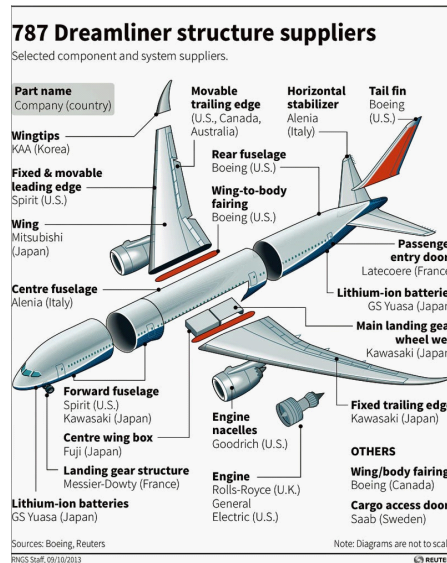


Figure 5: Reproduction of figure in “Boeing’s 787 Dreamliner Is Made of Parts from All over the World,” *Business Insider*, October 10, 2013.

The intention of Boeing’s organizational strategy was to speed up the development of the Dreamliner and save production costs. This is not what happened. As an article in Reuters reported at the time: “On a blustery and drizzly December afternoon in the Pacific Northwest, about 20 airplanes sat engineless and inert near the runway at a Boeing manufacturing plant...The program that produced these unfinished 787s is nearly three years behind schedule and, by some estimates, at least several billion dollars over budget.”⁵ The underlying reason for these delays and cost overruns were coordination problems among the suppliers and between them and Boeing. These problems proved so severe that Boeing was eventually forced to abandon its organizational strategy and bring the production of different modules back in-house: “Some of the parts arriving in Everett did not fit together, and late deliveries by producers of crucial sections of the plane

⁵ Kyle Peterson. “Special Report: A wing and a prayer: outsourcing at Boeing.” *Reuters*. January 20, 2011. In line with the description above, the article goes on to say: “The 787 is not merely a historic feat of engineering. The program also marks Boeing’s departure from its own time-honored manufacturing practices. Instead of drawing primarily from its traditional pool of aircraft engineers, mechanics and laborers that runs generations deep in the Puget Sound region around Seattle, Boeing leads an international team of suppliers and engineers from the United States, Japan, Italy, Australia, France and elsewhere, who make components that Boeing workers in the United States put together.”

stopped the entire assembly process...As a result, Boeing was forced to reverse some of its original outsourcing decisions; for example, in 2009 it spent \$1 billion in cash and credit to acquire its fuselage manufacturing partner Vought Aircraft Industries” (Brown and Garthwaite (2016), p.12). Even when products are highly modular, therefore, mirroring might fail because the need to coordinate across modules may necessitate intense communication between agents working on different ones.

In this section, we use the characterization of optimal communication networks in Proposition 2 to explore when the Mirroring Hypothesis does and does not hold. To allow for within-module communication, we relax the assumption that agents in the same module observe each others’ states and assume instead that each state is observed only by the associated agent. An organization then mirrors the production function if the principal places communication links within modules but not across.

DEFINITION. *An organization “mirrors” the production function if agent $i \in \mathcal{N}$ tells agent $j \in \mathcal{N}$ about his state if and only if they belong to the same module.*

For mirroring to be optimal, two conditions have to hold. First, profits have to be higher if each agent tells his state to the other agents in his own module, and no one else, than if he tells no agents at all. Second, conditional on each agent telling his state to the other agents in his own module, profits have to be higher if the agent refrains from telling any other agents in other modules. We can determine when these conditions are met by drawing on the results in the previous section. For the first condition, we know from (9) that if agent 1 tells his state to the other agents in his module, but no one else, his expected revenue is given by

$$R_1(\mathbf{C}_1(1)) = a_1^2 \sigma_1^2 \frac{1 - (n_1 - 2)p_1}{(1 + p_1)(1 - (n_1 - 1)p_1)},$$

which we can rewrite as

$$R_1(\mathbf{C}_1(1)) = a_1^2 \sigma_1^2 \left(1 + \frac{p_1^2 (n_1 - 1)}{(1 + p_1)(1 - (n_1 - 1)p_1)} \right).$$

The first term in brackets—the one—is the value of all walks from node 1 back to itself that only go through node 1 and the second is the value of all walks from node 1 back to itself that go through at least one other node in module \mathcal{M}_1 but do not go through any nodes in other modules. These latter walks are the ones that are created if agent 1 tells his state to the other agents in his own module but no one else. The next lemma then follows.

LEMMA 5. *Profits are higher if each agent tells his state to the other agents in his own module, and no one else, than if he tells his state to no agents at all if and only if*

$$a_i^2 \sigma_i^2 \frac{p_i^2}{(1 + p_i)(1 - (n_i - 1)p_i)} \geq \gamma \text{ for all } i \in \mathcal{N}. \quad (11)$$

The term on the left-hand side—the per node marginal revenue agent 1 generates when he tells the other agents in his module about his state—is increasing in n_i and p_i . For mirroring to be optimal, each module, therefore, has to have enough members, and coordination among them has to be sufficiently important.

At the same time, each module cannot have too many members, and coordination among them cannot be too important since otherwise it would be optimal to tell the agents in those modules about states in other modules. This property would violate the second condition for mirroring to be optimal: the absence of any communication across modules. Naturally, such across-module communication can never be optimal when the degree of coupling is sufficiently low. The more members modules have, though, and the more important it is to coordinate among them, the lower the degree of coupling needs to be for the absence of across-module communication to be optimal.

LEMMA 6. *Suppose each agent $i \in \mathcal{N}$ tells his state to the other agents in his module $\mathcal{M}_{m(i)}$. There exists a threshold degree of coupling $\bar{t}_i > 0$ such that profits are higher if agent i refrains from telling agents in other modules about his state if and only if $t \leq \bar{t}_i$. Adding modules to the production function decreases the threshold \bar{t}_i , as does increasing the module characteristics $n_{m'}$ or $p_{m'}$ for all $\mathcal{M}_{m'} \in \mathcal{M} \setminus \mathcal{M}_{m(i)}$.*

The conditions under which mirroring is optimal then follow directly.

PROPOSITION 3. *Mirroring is optimal if and only if (11) holds for all $i \in \mathcal{N}$ and $t \leq \min_{i \in \mathcal{N}} \bar{t}_i$.*

The proposition shows that for the Mirroring Hypothesis to hold, the conditions have to be just right. There cannot be too many modules and none of the modules can consist of too many decisions or require too much coordination, or else some agents should tell their information to agents in other modules. Arguably, this is why mirroring failed at Boeing. At the same time, there also cannot be modules that consist of too few members or require too little coordination, or else the agents in such modules should not even tell each other about their states. Mirroring, in other words, is associated with moderation. Modular production favors modular organization only if there is a limited number of modules, each of which consists of an intermediate number of decisions that require an intermediate degree of coordination.

A broader notion of the Mirroring Hypothesis allows for *modular-like* organizations, ones that contain clusters of modules whose agents communicate with each other but not with agents outside of the cluster. The management literature refers to such arrangements as “partial mirroring.”

DEFINITION. *An organization “partially mirrors” the production function if the set of modules \mathcal{M} can be partitioned into subsets such that (i.) agent $i \in \mathcal{N}$ tells agent $j \in \mathcal{N}$ about his state if*

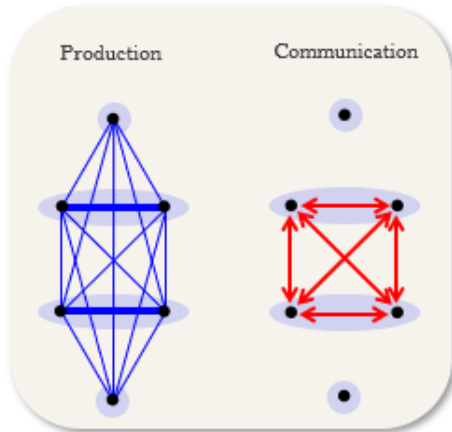


Figure 6: Illustration of an optimal communication network with a core-periphery structure. Left panel: production network consisting of four modules, two with two decisions and two with only one (where the blue shaded areas highlight the modules). Parameter values are as follows: need for coordination within the two-decision modules is 0.5, degree of coupling is 0.01, and the value of autonomous adaptation is one for all nodes. Right panel: the optimal communication network for any communication cost $\gamma \in (0.000434, 0.000801)$.

and only if their modules belong to the same subset and (ii.) there is at least one “cluster,” that is, a subset that contains two or more modules.

Partial mirroring can be optimal in our setting, and Figure 6 provides an example. If it is, though, it has to take a particular form.

PROPOSITION 4. *When partial mirroring is optimal, the organization contains one cluster of modules, and the modules that form the cluster are the most cohesive ones, that is, those $\mathcal{M}_m \in \mathcal{M}$ with the largest values of x_m .*

The result follows from the optimality of hierarchies we discussed in the previous section. If there were multiple clusters, hierarchies would not be nested, which cannot be optimal. Suppose, for instance, that one cluster consists of modules \mathcal{M}_1 and \mathcal{M}_2 and another of modules \mathcal{M}_3 and \mathcal{M}_4 . If it is optimal for an agent in \mathcal{M}_1 to tell his state to agents in \mathcal{M}_2 but not to those in \mathcal{M}_3 , then \mathcal{M}_2 has to be more cohesive than \mathcal{M}_3 . But if \mathcal{M}_2 is more cohesive than \mathcal{M}_3 , it cannot be optimal for an agent in \mathcal{M}_4 to tell his state to agents in \mathcal{M}_3 but not to those in \mathcal{M}_2 . In contrast, the existence of a single cluster is consistent with the optimal design of communication networks, provided it consists of the most cohesive modules. The optimal communication network then has a core-periphery structure in which the core consists of the most cohesive modules, whose agents all

communicate with each other, and the periphery consists of the less cohesive ones, whose agents only communicate with others in the same module.

This takes us back to Boeing and its response to the failure of its initial organizational strategy. By insourcing the production of some modules, such as the fuselage, while continuing to leave the production of others to its suppliers, Boeing created a core-periphery structure in which the inhouse modules formed the core and the outsourced ones the periphery. To the extent that the tail fin and the fuselage, as well as the other modules Boeing brought in-house, were the most cohesive ones, this response is consistent with the optimal design of communication networks in our model.

8 Robustness

Some of the assumptions in the model are critical for our results, while others are merely convenient. In the application in the previous section we already relaxed one of the convenient assumptions, the assumption that each agent observes all the states in his own module and does not have to be told about them. Beyond convenience, this assumption captures the notion that, because of physical proximity and shared expertise, agents working on the same module may know more about each others' local conditions than those working on different modules. Our results extend readily to an alternative specification in which each agent observes only his own state. We examine this alternative in Appendix B.

Some of the assumptions in the model are critical for the separability result in Proposition 1, without which an analytical characterization of optimal communication networks becomes much harder and, to us, intractable. One such critical assumption, and one that we share with Calvó-Armengol and de Martí (2008), Calvó-Armengol, de Martí, and Prat (2015), and Herskovic and Ramos (2020), is that agents do not re-transmit information they have received from others. If agent i talks to another agent, he tells him about his state θ_i but not about any other information he may have, such as the other states in his module. This assumption captures the notion that, while we model each agent's state as simply a number, it refers to a complex set of conditions and circumstances that only the associated agent can describe appropriately. If agents were able to re-transmit information, the separability result would fail. The communication links from agent i to other agents would then affect the overall cost, and thus optimal placement, of communication links from any other agent j .

The separability result also depends critically on the assumptions that states are independent and that communication is binary, that is, that agents learn the realization of a state either perfectly or not at all. In the absence of either assumption, the proof of Proposition 1 does not go through. We

share the assumption that states are independent with Calvó-Armengol, de Martí, and Prat (2015) and the binary communication structure with Dessein and Santos (2006) and Calvó-Armengol and de Martí (2008).

Finally, the separability result depends on the absence of incentive conflicts. As mentioned earlier, we share this assumption with the literature on team theory. To see how incentive conflicts affect the separability result, suppose that each agent only internalizes a fraction $\mu \in [0, 1]$ of the needs to coordinate and acts as if the production network were given by $\mu\mathbf{P}$ rather than \mathbf{P} (for instance because they put more weight on their own revenue or profits, as in Athey and Roberts (2001) and Alonso, Dessein, and Matouschek (2008)). The rest of the model is as in Section 3. The next proposition shows how such incentive conflicts affect the separability result.

PROPOSITION 5. *If agents internalize only a fraction $\mu \in [0, 1]$ of the needs to coordinate, an optimal communication network solves*

$$\max_{\mathcal{C}} \sum_{i=1}^N a_i \text{Cov}(d_i^*, \theta_i) + (1 - \mu) \sum_{i=1}^N \sum_{j=1}^N p_{ij} \text{Cov}(d_i^*, d_j^*) - \gamma \sum_{i=1}^N (\mathbf{C}_i \mathbf{1} - n_{m(i)}), \quad (12)$$

where

$$\text{Cov}(d_i^*, \theta_i) = a_i \sigma_i^2 \omega_{ii}(\mathbf{C}_i, \mu)$$

and

$$\text{Cov}(d_i^*, d_j^*) = \sum_{s=1}^N a_s^2 \sigma_s^2 \omega_{is}(\mathbf{C}_s, \mu) \omega_{js}(\mathbf{C}_s, \mu),$$

and where $\omega_{ij}(\mathbf{C}_j, \mu)$ denotes the ij th entry of $(\mathbf{I} - (\text{diag} \mathbf{C}_j) \mu \mathbf{P} (\text{diag} \mathbf{C}_j))^{-1}$.

The only new term in the principal's objective function (12) is the weighted sum of the covariances between each decision pair. Its presence implies that if agents are biased against coordination, it is no longer enough for the principal to ensure that each decision is sufficiently adapted to its state. Instead, she also needs to take into account how communication affects coordination and what she can do to ensure decisions co-vary more strongly with each other. The challenge this property poses is that the extent to which two decisions co-vary with each other depends on the overlap of the decision makers' information sets, that is, on which states they are both informed about. The principal can, therefore, no longer consider each agent in isolation and ask whom he should tell about his state. She has to consider all agents at once and take into account how communication links from one agent affect the optimal location of such links from the others. Since the objective function is still supermodular, the principal can still use standard algorithms to solve for optimal communication networks in polynomial time. Finding an analytical solution, however, is now more challenging.

The key assumption in the entire paper is that the production function has a non-overlapping community structure. This assumption allows us to capture the notion that products are modular, which are the type of products we are interested in. To generalize this structure, one could allow for different degrees of coupling—different ts —for decisions in different module pairs. Since a module may consist of a single decision, though, such a production function would constitute a general, unweighted network with little structure to base a characterization on. To see what can still be said in this case, suppose that the production network \mathbf{P} can take any form, provided it still satisfies $p_{ii} = 0$, $p_{ij} = p_{ji}$, and $\sum_{j=1}^N p_{ij} < 1$ for all $i, j \in \mathcal{N}$. The separability result in Proposition 1, and the lemmas that precede it, continue to hold for such more general production functions. As such, the principal can still determine the optimal communication network by considering each agent in isolation. Moreover, the principal’s objective is still supermodular and can, therefore, be maximized using standard algorithms.

What can no longer hold is the characterization of optimal communication networks in Proposition 2, which is specific to the non-overlapping community structure. Optimal communication can now take many forms and need not give rise to hierarchies. The specific form it takes depends on the specific structure of the production network. There are, however, some properties of optimal communication networks that hold across production networks.

PROPOSITION 6. *As long as the production network \mathbf{P} satisfies $p_{ii} = 0$, $p_{ij} = p_{ji}$, and $\sum_{j=1}^N p_{ij} < 1$ for all $i, j \in N$, optimal communication networks \mathbf{C}_i^* are increasing in the value of autonomous adaptation $a_i^2 \sigma_i^2$ and the needs for coordination p_{ij} , and decreasing in communication costs γ .*

The proposition shows that, independent of the specific structure of the production network, the principal will only ever respond to an increase in the value of adaptation or the need for coordination, or a decrease in the cost of communication, by adding communication links. These comparative statics hold because the principal’s objective function is supermodular and has either increasing or decreasing differences in the various parameters (Topkis 1978, Milgrom and Shannon 1994).

9 Conclusions

Since the middle of the last century, modular production has emerged as a prevalent form of production. The rise of modular production has been widely observed and documented and has been explored extensively in management and computer science. The goal of this paper is to take a first step towards understanding the economic implications of the rise of modular production.

As a first step, we focused on the immediate implications of modular production for the internal organization of firms and abstracted from any broader implications for their boundaries and the structure of industries. Even in this narrow context, many open questions remain. An important practical issue we put aside is the role of “interfaces” which ensure that different modules fit with each other. One way to think about such interfaces in our model is as a limited set of decisions that are made and announced before agents make the remaining decisions.

Another widely-discussed issue we did not address is “parallel processing,” the notion that modular production allows firms to accelerate production by having different agents work on different modules simultaneously (Parnas 1972). One way to get at this issue in our model is to suppose that the principal can hire agents and decide which decisions each agent is in charge of. Each agent first spends time learning the states associated with his decisions, taking one period per state to do so. After all the agents have learned their states, they make their decisions simultaneously and without spending any further time on communication. A patient principal would hire a single agent and have him make the first-best decisions after N periods but an impatient principal may prefer to hire M agents, put each in charge of one module, and have them make second-best decisions sooner. We leave the investigation of both interfaces and parallel processing, as well as other issues related to internal organization, for future research.

The impact of modular production on the economy is unlikely to be confined to changes in the internal organization of firms. Baldwin and Clark (1997), for instance, observe that while the introduction of the System/360 did lead to immediate changes in IBM’s internal organization, its more enduring impact was to cause entry into the computer industry in the following decades. The entrants were often small, entrepreneurial firms who focused on the development and production of individual modules and whose innovative products allowed them to compete successfully with IBM’s own, inhouse module makers. In this telling, the introduction of the System/360 in the 1960s sowed the seeds for the subsequent disintegration of IBM and the other large mainframe manufacturers and gave rise to the competitive and innovative computer industry of today.⁶ There are many reasons why modular production may affect the boundaries of firms and the structure and inventiveness of industries and we leave their exploration for future research.

A question that goes beyond the impact of modular production on economic activity is what

⁶As Baldwin and Clark (1997) observe: “But modularity also undermined IBM’s dominance in the long run, as new companies produced their own so-called plug-compatible modules—printers, terminals, memory, software, and eventually even the central processing units themselves—that were compatible with, and could plug right into, the IBM machines. By following IBM’s design rules but specializing in a particular area, an upstart company could often produce a module that was better than the ones IBM was making internally. Ultimately, the dynamic, innovative industry that has grown up around these modules developed entirely new kinds of computer systems that have taken away most of the mainframe’s market share.”

explains its rise in the first place. Herbert Simon argued that modularity facilitates adaptation by confining adaptive changes to individual modules within a system. The argument that modularity allows parallel processing provides another reason why it may have adaptive advantages. In line with these intuitions, firms such as IBM explain their development of modular products with the need to adapt quickly to the changing capabilities of their suppliers and needs of their customers. Yet, a full explanation for the rise of modular production also needs to account for its costs. It may be easier to adapt a modular product to its environment but, for a given environment, one would expect limitations in across-module interactions to affect its quality. After all, products have not always been modular, and even today many are not, suggesting that such designs also have significant downsides. Answering the questions of when and why firms develop modular products, and what trade-offs they face when they are doing so, would require moving beyond one of the foundational economic modeling assumptions, that production functions are given by nature and not designed by firms. As such, it is the most challenging question this paper highlights and, like the other open questions we sketched above, we leave it for future research.

Finally, this paper also raises empirical issues. Newly emerging data sets contain detailed information about communication between employees of real-world firms (Impink, Prat, and Sadun (2021) and Yang et al. (2021)). Our model makes specific predictions about the pattern of such communication in firms that make modular products. In particular, the prediction that optimal communication has a nested, hierarchical structure has a number of implications that are, at least in principle, observable, such as the emergence of core-periphery structures and the absence of multiple cores or clusters. While testing the model is naturally difficult, we hope that this paper provides some stimulation and direction to the budding empirical literature on within-firm communication.

References

- [1] ALONSO, RICARDO, WOUTER DESSEIN, AND NIKO MATOUSCHEK. 2008. When Does Coordination Require Centralization? *American Economic Review*. 98(1): 145-79.
- [2] ARROW, KENNETH. 1974. *The Limits of Organization*. W.W. Norton & Company.
- [3] ATHEY, SUSAN AND JOHN ROBERTS. 2001. Organizational Design: Decision Rights and Incentive Contracts. *American Economic Review*. 91(2): 200-5.
- [4] BALA, VENKATESH AND SANJEEV GOYAL. 2000. A Noncooperative Model of Network Formation. *Econometrica* 68(5): 1181-229.
- [5] BALDWIN, CARLISS AND KIM CLARK. 1997. Managing in an Age of Modularity. *Harvard Business Review*: 84-93.
- [6] ——— AND ——— . 2000. *Design Rules: The Power of Modularity (Vol.1)*. M.I.T. Press.
- [7] BERGEMANN, DIRK, TIBOR HEUMANN, AND STEPHEN MORRIS. 2017. Information and Interaction. Cowles Foundation Discussion Paper No. 2088.
- [8] BROWN, JENNIFER AND CRAIG GARTHWAITE. 2016. Global Aircraft Manufacturing, 2002-2011. *Kellogg School of Management Case Study*. KEL 938.
- [9] COLFER, LYRA AND CARLISS BALDWIN. 2016. The Mirroring Hypothesis: Theory, Evidence, and Exceptions. *Industrial and Corporate Change*. 25(5): 709-38.
- [10] CONWAY, MELVIN. 1968. How Do Committees Invent? *Datamation*. 14(4): 84-93.
- [11] DESSEIN, WOUTER, ANDREA GALEOTTI, AND TANO SANTO., 2016. Rational Inattention and Organizational Focus. *American Economic Review*. 106(6): 1522-36.
- [12] FORTUNATO, SANTO. 2010. Community Detection in Graphs. *Physics Reports*. 486(3-5): 75-174.
- [13] GARICANO, LUIS. 2000. Hierarchies and the Organization of Knowledge in Production. *Journal of Political Economy*, 108(5): 874-904.
- [14] GARUD, RAGHU, ARUN KUMARASWAMY, AND RICHARD LANGLOIS, eds. 2009. *Managing in the Modular Age: Architectures, Networks, and Organizations*. John Wiley & Sons.

- [15] GUIMERA, ROGER, STEFANO MOSSA, ADRIAN TURTSCHI, AND LA NUNES AMARAL. 2005. The Worldwide Air Transportation Network: Anomalous Centrality, Community Structure, and Cities' Global Roles. *Proceedings of the National Academy of Sciences*. 102(22): 7794-99.
- [16] HERSKOVIC, BERNARD AND JOAO RAMOS. 2020. Acquiring Information through Peers. *American Economic Review*. 110(7): 2128-52.
- [17] IMPINK, STEPHEN, ANDREA PRAT, AND RAFFAELLA SADUN. 2021. Communication within Firms: Evidence from CEO Turnovers. NBER Working Paper Series, No. 29042.
- [18] LANGLOIS, RICHARD AND PAUL ROBERTSON. 1992. Networks and Innovation in a Modular System: Lessons from the Microcomputer and Stereo Component Industries. *Research Policy*. 21(4): 297-313.
- [19] MEUNIER, DAVID, RENAUD LAMBIOTTE, ALEX FORNITO, KAREN ERSCHKE, AND EDWARD BULLMORE. Hierarchical Modularity in Human Brain Functional Networks. *Frontiers in Neuroinformatics*. 3: 37.
- [20] MCCORD, KENT AND STEVEN EPPINGER. 1993. Managing the Integration Problem in Concurrent Engineering. M.I.T. Sloan School of Management, Cambridge, MA, Working Paper No. 3594.
- [21] MILGROM, PAUL AND CHRIS SHANNON. 1994. Monotone Comparative Statics. *Econometrica*. 157-80.
- [22] MORRIS, STEPHEN. 2000. Contagion. *The Review of Economic Studies*. 67(1): 57-78.
- [23] MUROTA, KAZUO. 2003. Discrete Convex Analysis. Society for Industrial and Applied Mathematics.
- [24] PETERSON, KYLE. 2011. Special Report: A wing and a prayer: outsourcing at Boeing. *Reuters*. January 20, 2011.
- [25] SIMON, HERBERT. 1962. The Architecture of Complexity. *Proceedings of the American Philosophical Society*. 106 (6): 467-482.
- [26] TADELIS, STEVEN AND OLIVER WILLIAMSON. 2012. 4. Transaction Cost Economics. Princeton University Press: 159-90.

- [27] THOMPSON, JAMES. 1967. *Organizations in Action: Social Science Bases of Administrative Theory*. McGraw-Hill.
- [28] TOPKIS, DONALD. 1978. Minimizing a Submodular Function on a Lattice. *Operations Research*. 26(2): 305-21.
- [29] LONGQI YANG, DAVID HOLTZ, SONIA JAFFE, SIDDHARTH SURI, SHILPI SINHA, JEFFREY WESTON, CONNOR JOYCE, NEHA SHAH, KEVIN SHERMAN, BRENT HECHT AND JAIME TEEVAN. 2021. The Effects of Remote Work on Collaboration Among Information Workers. *Nature Human Behaviour*: 1-12.