# The Market For Data Privacy

Tarun Ramadorai

Antoine Uettwiller

Ansgar Walther

# Data Privacy in the Internet Era

Firms collect, share and aggregate data about a wide range of consumers' online and offline activities

Varian, 2009; Krishnamurthy and Wills, 2009; FTC, 2014

Economics principles are subtle:

► Classical: Consumer data improves efficiency of allocations

Stigler, 1980; Posner, 1981; Goldfarb and Tucker, 2011

► Second best: Concerns about insurance, price discrimination, negative externalities

Hirshleifer, 1971; Taylor, 2004; Varian, 2009

**How does the market for data privacy operate?**

# The Market for Data Privacy

Demand: Many consumers are passive, "consent fatigue"

Goldfarb and Tucker; 2012; Acquisti et al., 2015; Campbell et al., 2018

► Privacy paradox: stated preferences vs. behavior and WTP

► Reassurance by mere presence of legal text

   Norberg et al., 2007; Acquisti, 2016; Athey et al., 2017

Understanding *supply of privacy* is important in this context

**This paper: What determines firms' privacy contracts and data sharing policies?**

# This Paper

**Data collection:** For a comprehensive set of US firms, we measure

1. What they say: Privacy policy text
2. What it means: Evaluation of these policies by a legal expert
3. What they do: Third party cookies on websites
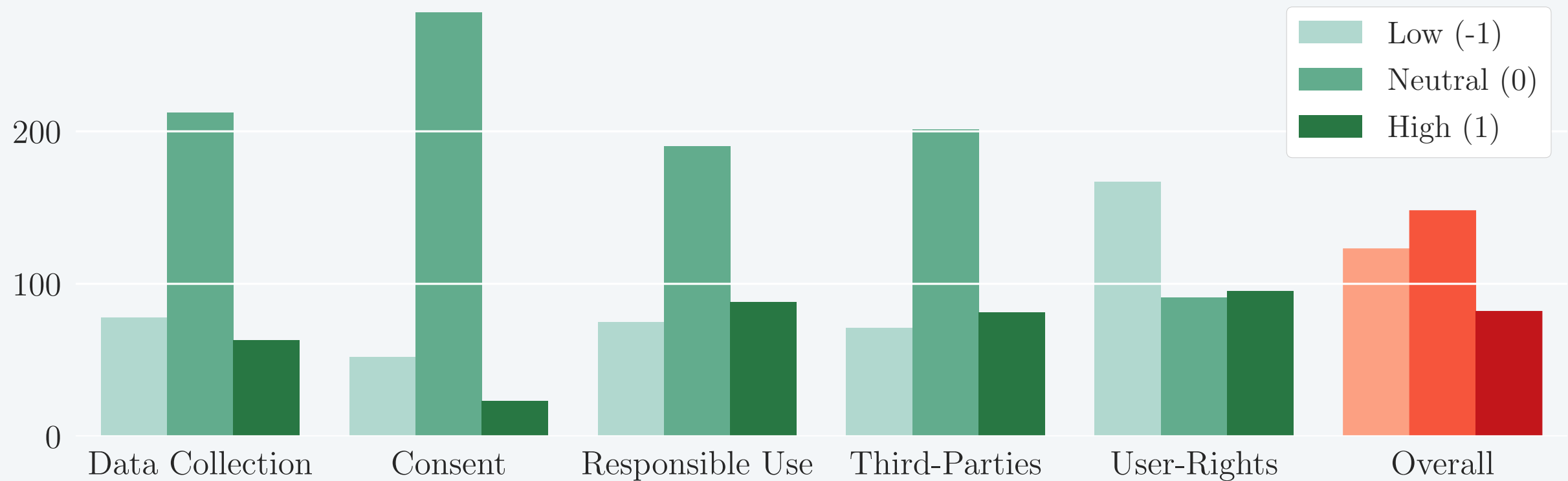
**Stylized facts using variation across firms:**

► No standard industry-level boilerplate
► Detailed policies are **associated with more sharing** (fig leaves?)
► Systematic variation across firm characteristics

  ► Size and technical sophistication

**Theory:** Determinants of firms' data sharing and privacy policies

# Data

# Privacy Policies

$N =$5377 firms in Compustat US

Finding privacy policies:

► Automated google search

► Web crawling

► Manual checking

Visibility: "Privacy" link on website

## Access and Visibility



## Word Cloud

# Expert Evaluation
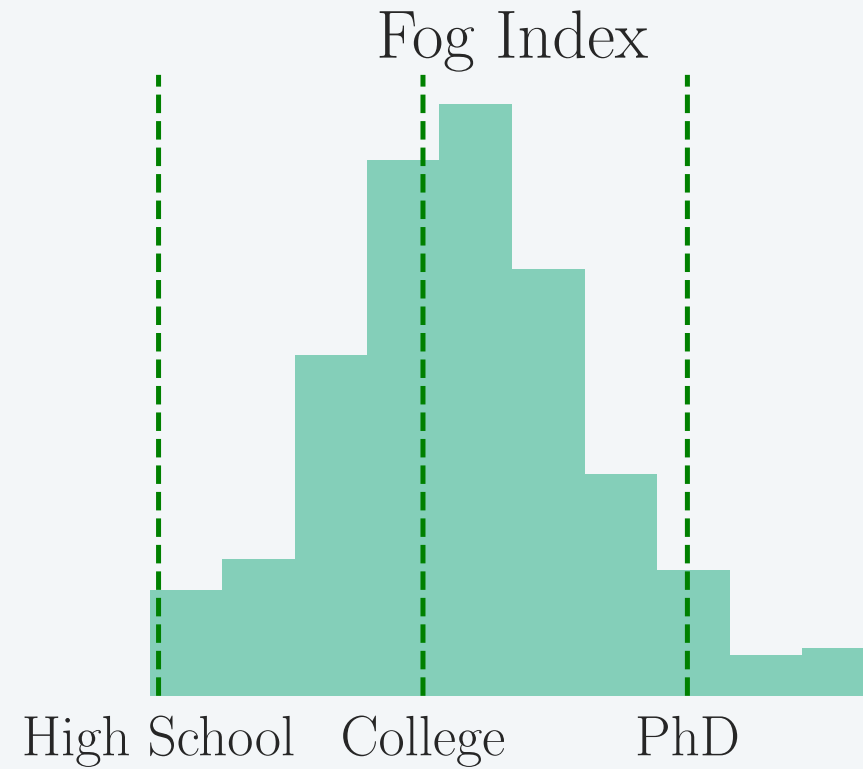
## 10% Sample of Policies



$$Legal\ Clarity\ Index = \text{Frequency of top 100 "High Overall" bigrams}$$

$$- \text{Frequency of top 100 "Low Overall" bigrams}$$

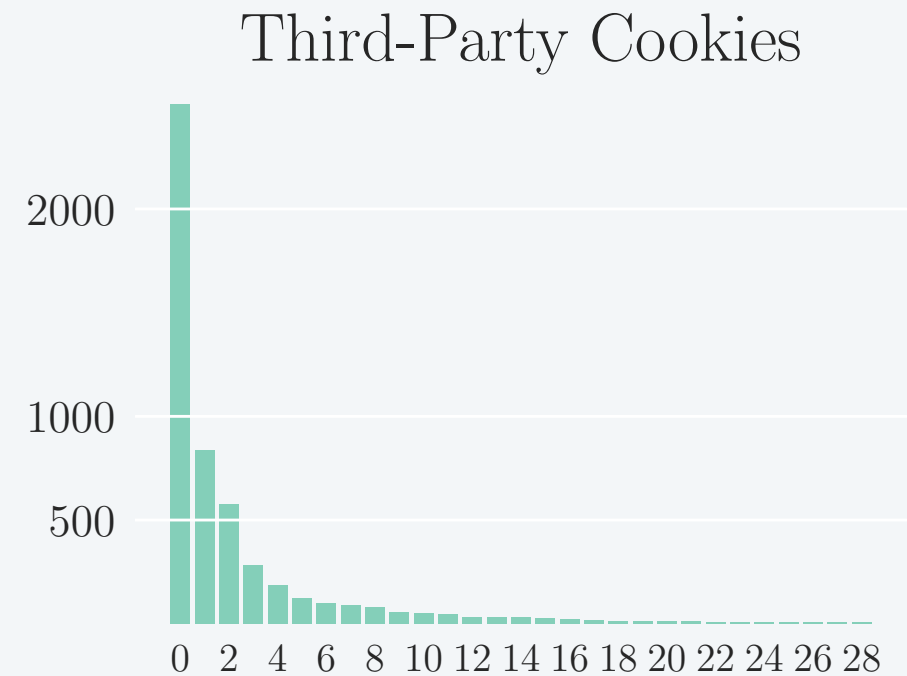# Readability and Third Party Data Sharing

"Fog" readability index: Years of formal education needed to read a document

Gunning, 1952

### Fog Index

High School    College    PhD

OpenWPM scraper counts cookies on firm's website

Englehardt and Narayanan, 2016

### Third-Party Cookies

2000

1000

500

0 2 4 6 8 10 12 14 16 18 20 22 24 26 28
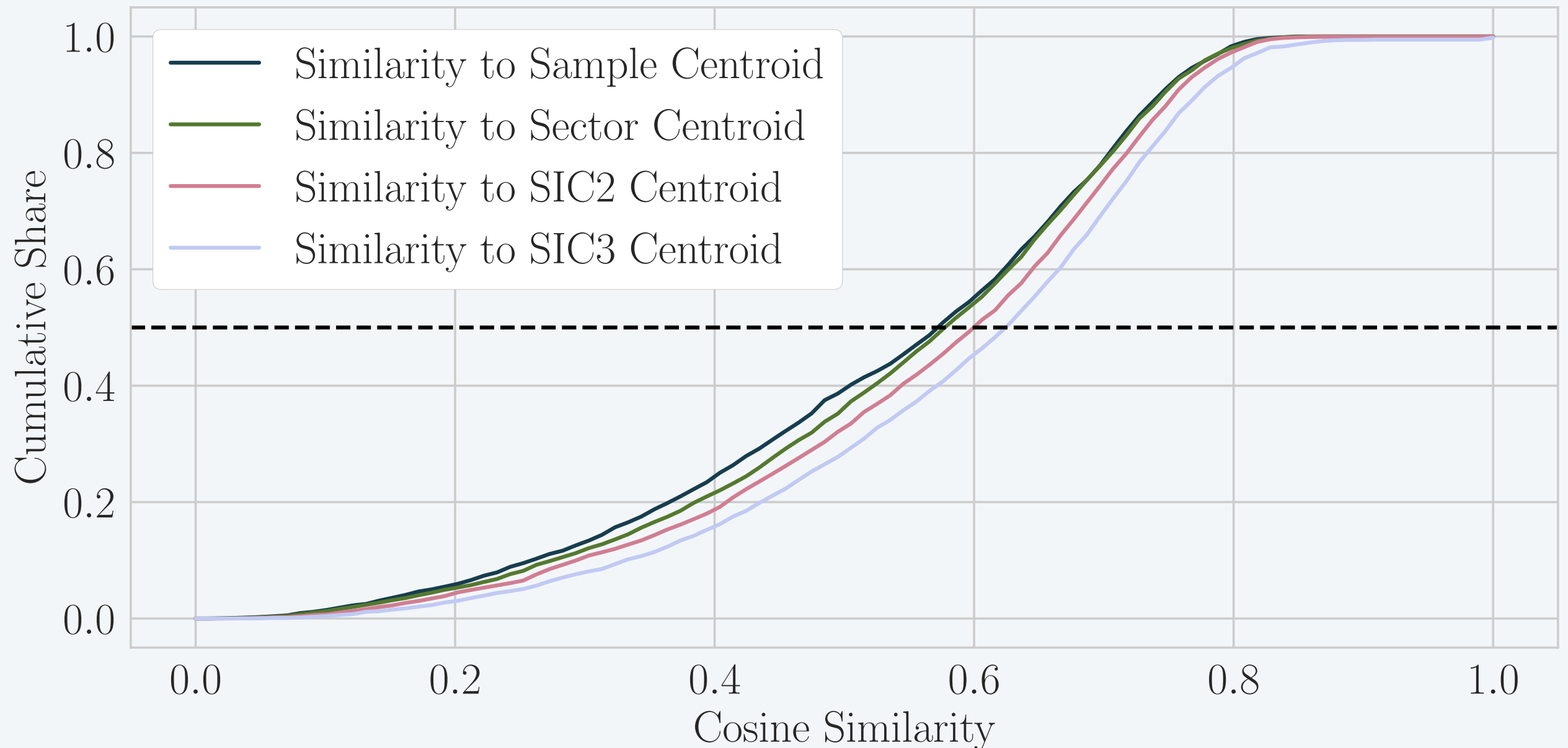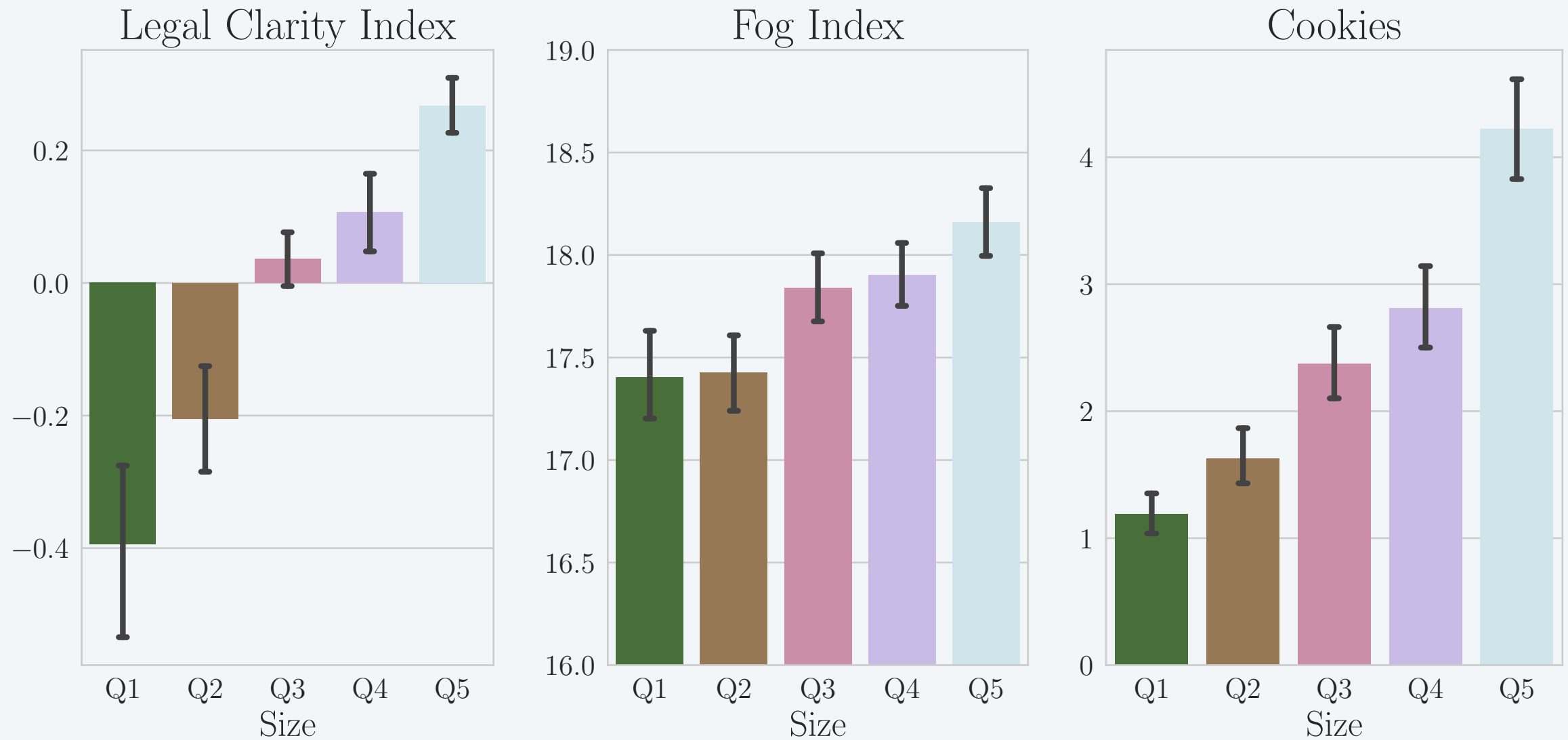
# Stylized Facts

# Variation: No Industry Boilerplates

## Similarity of Word Frequency Vectors Across Policies



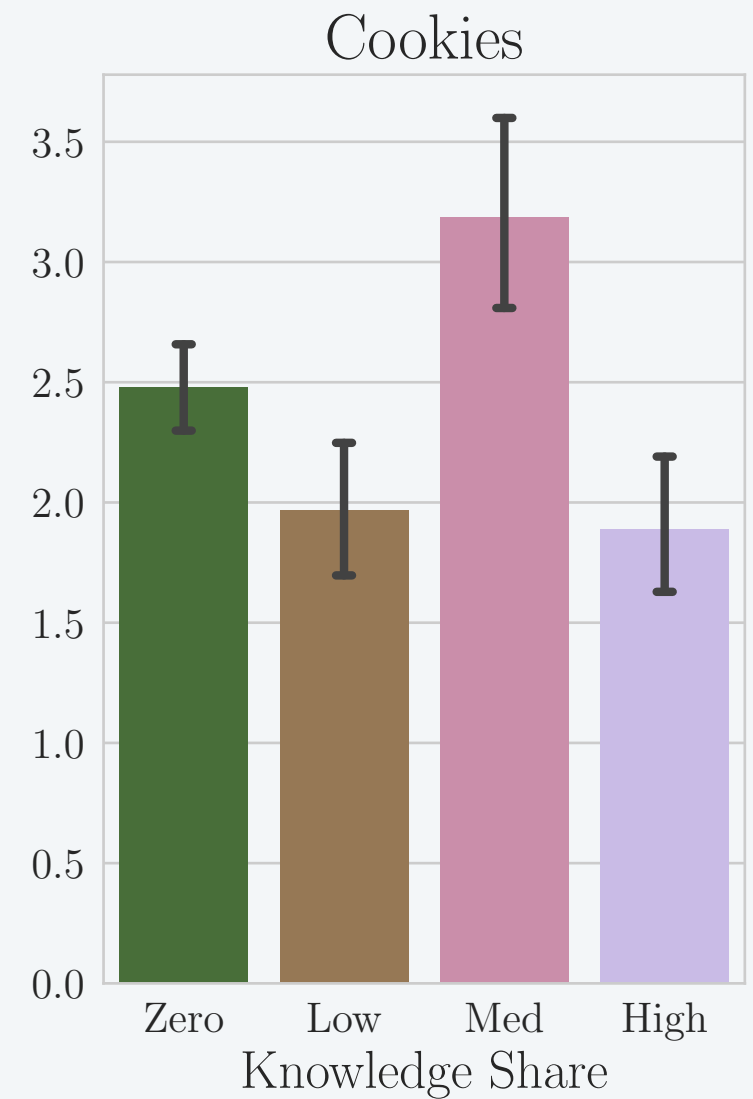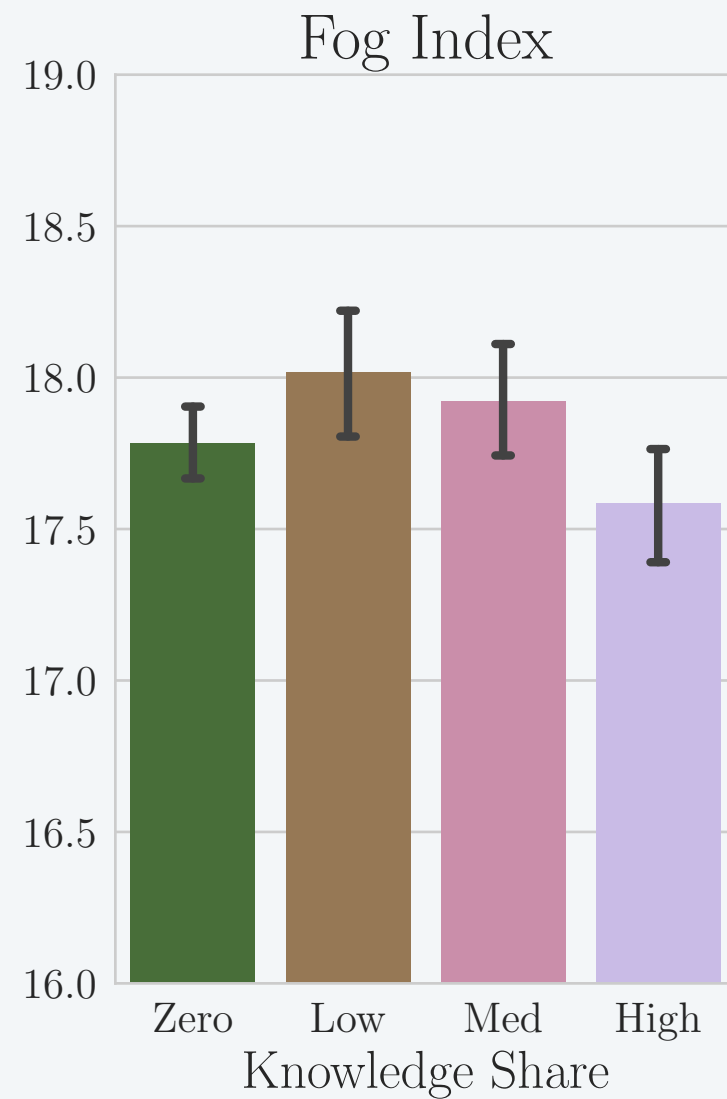Only slight increase in similarity in latent "topic" space
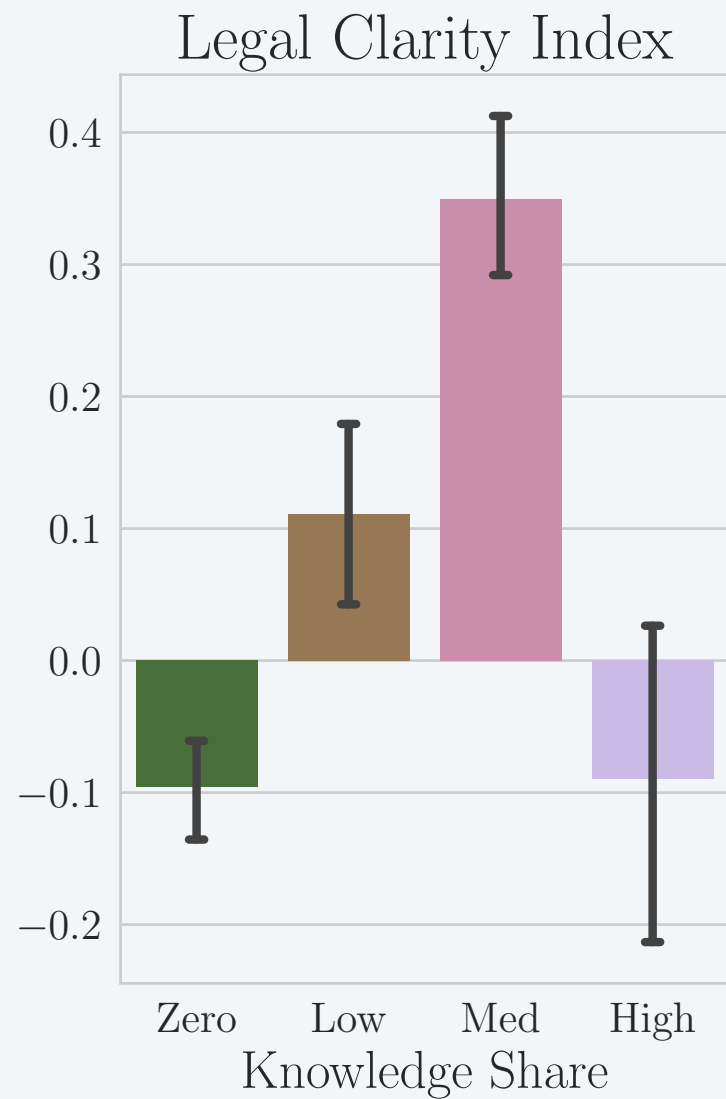
# Firm Size, Policies and Behavior



Large firms also have longer policies which are easier to find

# Knowledge Share, Policies and Behavior

## Capital Accumulated through R&D / Total Assets
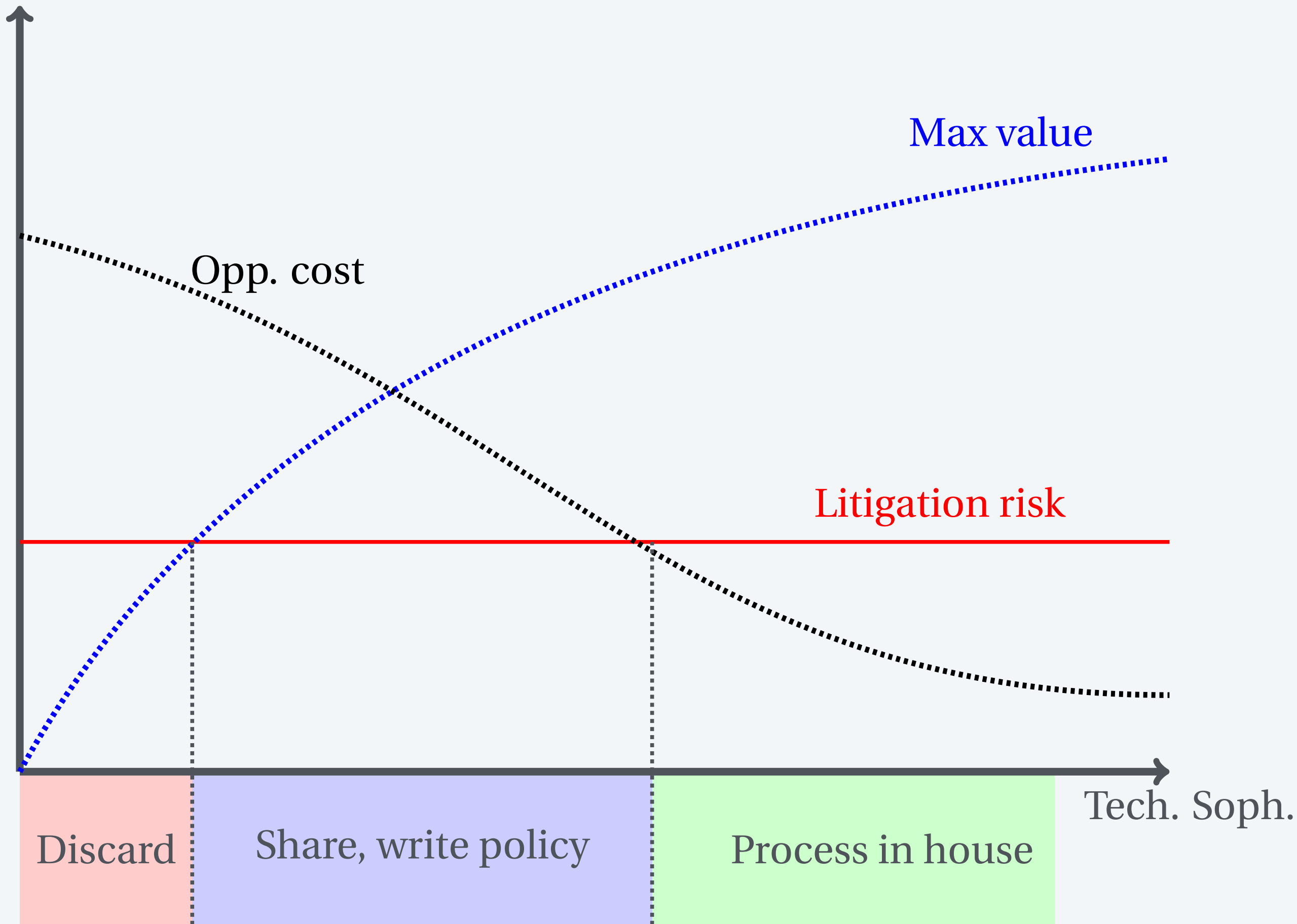
# Theory of Data Sharing

# Model

**Firm** can monetize its data by:

1. Processing in-house

2. Sharing with a specialist **data intermediary**

Sufficient statistics:

► Max value of data

► Opportunity cost of in-house processing

► Litigation risk (alleviated by costly privacy policy)

Max value

Opp. cost

Litigation risk

Tech. Soph.

Discard

Share, write policy

Process in house

# Partial Effects of Knowledge Capital

Controlling for firm size, market share, industry FE:

| | (1)<br>Policy Found | (2)<br>Policy Visible | (3)<br>Log Words | (4)<br>Overall Score | (5)<br>Fog Index | (6)<br>$3^{rd}$-Party Trackers |
|---|---|---|---|---|---|---|
| Log Market Value | 0.0421*** | 0.0484*** | -0.00597 | 0.0426*** | 0.0296 | 0.330*** |
| | (12.22) | (12.13) | (-0.61) | (4.44) | (1.14) | (8.20) |
| Knowledge Share | 0.847*** | 0.695*** | 2.405*** | 2.605*** | 0.501 | 4.447*** |
| | (8.33) | (5.89) | (8.80) | (9.78) | (0.69) | (3.76) |
| Knowledge Share$^2$ | -0.813*** | -0.793*** | -2.821*** | -3.811*** | -0.264 | -7.114*** |
| | (-4.90) | (-4.12) | (-6.30) | (-8.74) | (-0.22) | (-3.69) |
| Log Market Share | 0.0157*** | -0.0105*** | 0.0874*** | 0.0615*** | 0.100*** | 0.119*** |
| | (5.41) | (-3.11) | (10.49) | (7.57) | (4.54) | (3.52) |
| Observations | 5140 | 5140 | 3918 | 3918 | 3918 | 4951 |

# Conclusions

▶ We assemble comprehensive data for studying the market for privacy, focusing on the supply side

▶ Stylized facts on cross-firm variation

    ▶ Clear policies $\Rightarrow$ more sharing

▶ Simple testable theory of data sharing

# Public Resources

**www.github.com/ansgarw/privacy**

► All our data for work with Compustat US firms

► Python code, demos and documentation

► Get policies and their attributes for *any* sample of firms or websites

## Simplest Example

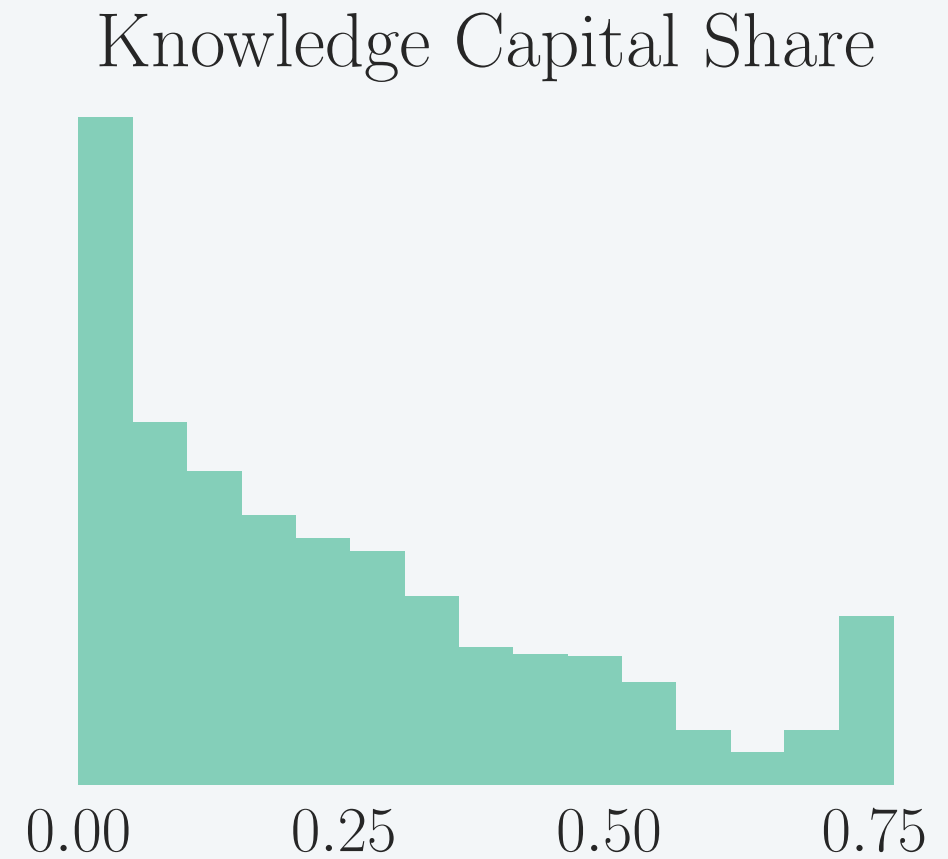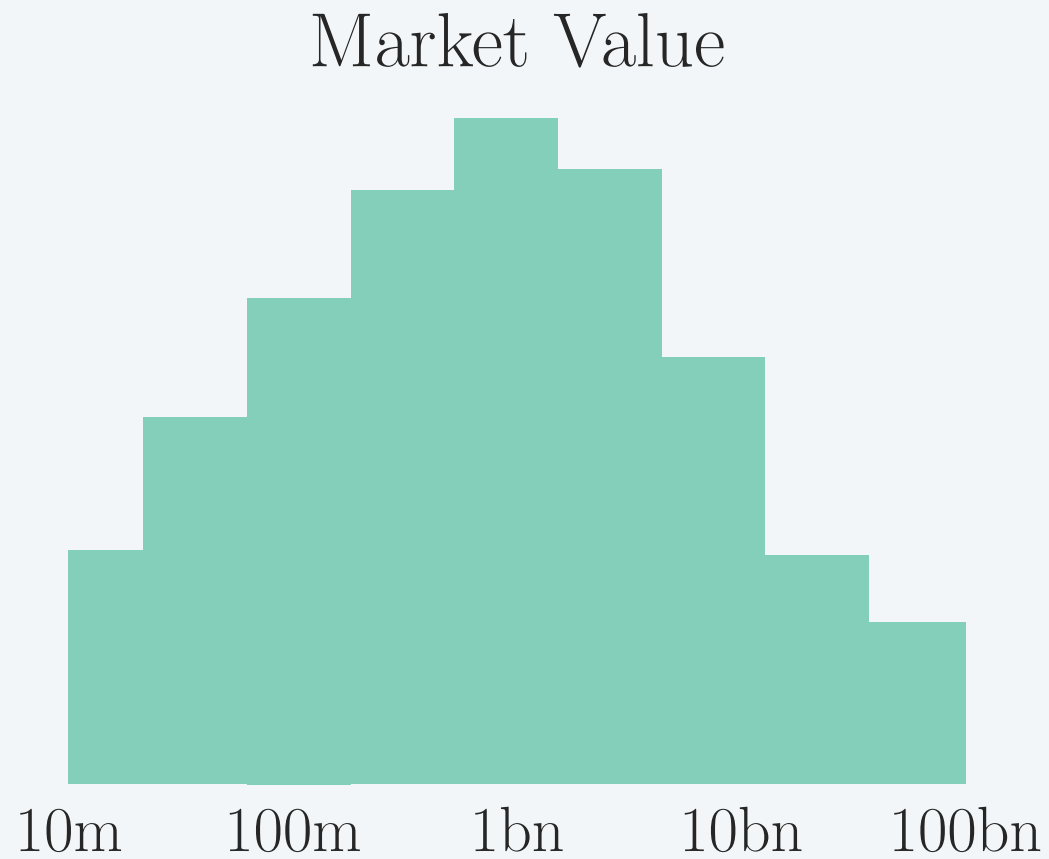Here are 5 lines of code that find the policy for American Airlines:

```
from src.urls import crawlPrivacy, filterPrivacy
from src.text import findPolicy
status, urls = crawlPrivacy('www.aa.com',clicks=2) # crawls candidate URLs
ranked = filterPrivacy(sum(urls,[])) # filter and rank by likelihood of being privacy policy
status, policy, url = findPolicy(ranked) # scrape highest ranked page that contains 'privacy'
```



```
Legal clarity of www.aa.com: 1.136
Legal clarity of www.ba.com: 1.691
```

# Firm Characteristics

Market Value

Knowledge Capital Share

10m   100m   1bn   10bn   100bn

0.00   0.25   0.50   0.75

$$\text{Knowledge Share} = \frac{\text{Capital accumulated through R\&D}}{\text{Total Assets}}$$

Peters and Taylor, 2017   Back to sorts

# Legal Clarity Index



High and low score policies look different, so we construct:

*Legal Clarity Index* = Frequency of top 100 "High" bigrams

− Frequency of top 100 "Low" bigrams

Similar results with an index that uses supervised machine learning