# Data-intensive Innovation and the State: Evidence from AI Firms in China

Martin Beraja    David Yang    Noam Yuchtman
MIT              Harvard       LSE

ABFER webinar
November 3, 2021

# Motivation: government data as input in AI innovation

▶ AI innovation is **data-intensive**

    ▶ Many recent AI advances made with decades-old algorithms applied to newly available big data

# Motivation: government data as input in AI innovation

- ▶ AI innovation is **data-intensive**
  - ▶ Many recent AI advances made with decades-old algorithms applied to newly available big data

- ▶ Literature has focused on how data collected by **private** firms shapes AI innovation (Agrawal et al., 2019; Jones and Tonetti, 2020)

# Motivation: government data as input in AI innovation

- ▶ AI innovation is **data-intensive**
  - ▶ Many recent AI advances made with decades-old algorithms applied to newly available big data

- ▶ Literature has focused on how data collected by **private** firms shapes AI innovation (Agrawal et al., 2019; Jones and Tonetti, 2020)

- ▶ Yet, throughout history, **states** have also collected massive quantities of data (Scott, 1998)

- ▶ The state has a large role in many areas
  - ▶ Public security, health care, education, basic science...

# Motivation: government data as input in AI innovation

- ▶ AI innovation is **data-intensive**
    - ▶ Many recent AI advances made with decades-old algorithms applied to newly available big data

- ▶ Literature has focused on how data collected by **private** firms shapes AI innovation (Agrawal et al., 2019; Jones and Tonetti, 2020)

- ▶ Yet, throughout history, **states** have also collected massive quantities of data (Scott, 1998)

- ▶ The state has a large role in many areas
    - ▶ Public security, health care, education, basic science...

    $\implies$ **Government data** can exceed privately-collected data in magnitude/scope; or lack good substitutes altogether

# Motivation: China's facial recognition AI sector

▶ A common way in which AI firms **gain access** to valuable government data is by **providing services** to the state

# Motivation: China's facial recognition AI sector

► A common way in which AI firms **gain access** to valuable government data is by **providing services** to the state

► Think about **facial recognition AI firms in China**...

   ► Train algorithms with, e.g., video streams of faces from many angles

   ► The state's public security units collect this form of data through their surveillance apparatus, and contract AI firms for services

   ► AI firms gaining access to surveillance data can use it to train algorithms and develop software

# This paper

Does access to **government data** when providing AI services
to the state stimulate **commercial** AI innovation?

# This paper

Does access to **government data** when providing AI services to the state stimulate **commercial** AI innovation?

**The mechanism(s)**

1. If gov't data and algorithms are **sharable** across uses, they can be used to develop AI products for commercial markets
   (e.g., a facial recognition platform for retail stores)

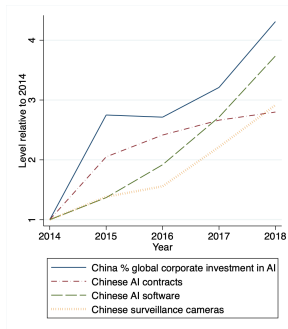2. Firms may **learn** to manage and utilize large datasets too

$\implies$ a procurement contract with access to gov't data can fuel commercial innovation, overcoming **crowd-out** from the contract

# This paper

Does access to **government data** when providing AI services to the state stimulate **commercial** AI innovation?

**The mechanism(s)**

1. If gov't data and algorithms are **sharable** across uses, they can be used to develop AI products for commercial markets

   (e.g., a facial recognition platform for retail stores)

2. Firms may **learn** to manage and utilize large datasets too

$\implies$ a procurement contract with access to gov't data can fuel commercial innovation, overcoming **crowd-out** from the contract

Evidence of this in China's facial recognition AI sector

# Two implications

1. Access to gov't data contributed to Chinese firms' emergence as leading innovators in facial recognition AI

   ▶ Indeed, this has coincided with the expansion of the government's procurement of AI and surveillance capacity

# Two implications

1. Access to gov't data contributed to Chinese firms' emergence as leading innovators in facial recognition AI

   ▶ Indeed, this has coincided with the expansion of the government's procurement of AI and surveillance capacity

2. Novel role for the state in data-intensive economies

   ▶ So far, emphasis on the regulation of privately-collected data due to antitrust or privacy concerns (Tirole, 2020; Aridor et al., 2020)

   ▶ AI procurement and policies of gov't data collection and provision could, **whether intentionally or not**, stimulate and shape the direction of innovation in a range of sectors

# Empirical challenges

Would like to compare software output changes after receipt of gov't procurement contracts giving access to more v. less data

# Empirical challenges

Would like to compare software output changes after receipt of gov't procurement contracts giving access to more v. less data

**Data challenges**

1. Dataset linking AI firms to govt. contracts did not exist

2. Dataset on AI firms' software did not exist (our measure of *product innovation*). Also, critical for us to classify by use (commercial or not)

3. No available direct measures of firm-level use of gov't data

# Empirical challenges

Would like to compare software output changes after receipt of gov't procurement contracts giving access to more v. less data

**Data challenges**

1. Dataset linking AI firms to govt. contracts did not exist

2. Dataset on AI firms' software did not exist (our measure of *product innovation*). Also, critical for us to classify by use (commercial or not)

3. No available direct measures of firm-level use of gov't data

**Identification challenges**

1. Non-random assignment of gov't contracts

2. Contracts work through other mechanisms unrelated to data

# Data 1: linking AI firms to govt. contracts

1. Identify all facial recognition AI **firms**

   - 7,837 firms
   - Two sources: Tianyancha (People's Bank of China) and PitchBook (Morningstar)
   - Include: *(i)* firms specialized in facial recognition AI (e.g., Yitu); *(ii)* hardware firms that devote substantial resources to develop AI software (e.g., Hik-Vision); *(iii)* facial recognition AI units of large tech conglomerates (e.g., Baidu AI)

# Data 1: linking AI firms to govt. contracts

1. Identify all facial recognition AI **firms**

    - 7,837 firms
    - Two sources: Tianyancha (People's Bank of China) and PitchBook (Morningstar)
    - Include: *(i)* firms specialized in facial recognition AI (e.g., Yitu); *(ii)* hardware firms that devote substantial resources to develop AI software (e.g., Hik-Vision); *(iii)* facial recognition AI units of large tech conglomerates (e.g., Baidu AI)

2. Obtain universe of **government contracts**

    - 2,997,105 contracts
    - Source: Chinese Govt. Procurement Database (Ministry of Finance)

# Data 1: linking AI firms to govt. contracts

1. Identify all facial recognition AI **firms**

   - 7,837 firms
   - Two sources: Tianyancha (People's Bank of China) and PitchBook (Morningstar)
   - Include: *(i)* firms specialized in facial recognition AI (e.g., Yitu); *(ii)* hardware firms that devote substantial resources to develop AI software (e.g., Hik-Vision); *(iii)* facial recognition AI units of large tech conglomerates (e.g., Baidu AI)

2. Obtain universe of **government contracts**

   - 2,997,105 contracts
   - Source: Chinese Govt. Procurement Database (Ministry of Finance)

3. Link government **buyers** to AI **suppliers**

# Data 2: AI firms' software production

Registered with Min. of Industry and Information Technology
- Validation exercise: check against IPO Prospectus of MegVii

# Data 2: AI firms' software production

Registered with Min. of Industry and Information Technology
- Validation exercise: check against IPO Prospectus of MegVii

**Categorize by intended customers:**

1. **Commercial:** e.g., *visual recognition system for smart retail*;
2. **Government:** e.g., *smart city — real time monitoring system on main traffic routes*;
3. General: e.g., *a synchronization method for multi-view cameras based on FPGA chips*.

# Categorization: analyze text using machine learning

▶ Recurrent Neural Network (RNN) model using tensorflow

  - Corpus: 13,000 manually labeled software programs
  - Word-embedding: converted sentences to vectors based on word frequencies and used the words from full datasets as dictionary
  - Long Short-Term Memory (LSTM) algorithm: 2 layers of 32 nodes
  - 90% of corpus for training, 10% for validating
  - 10,000 training cycles are run for gradient descent on loss function

▶ Results robust to perturbing parameters of learning model

# Data 3: measuring access to government data

**Within AI public security contracts:** variation in the data collection capacity of the public security agency's local surveillance network

1. Identify non-AI contracts: police department purchases of street cameras
2. Measure quantity of advanced cameras in a prefecture at a given time
3. Categorize public security contracts as coming from "high" or "low" camera capacity prefectures
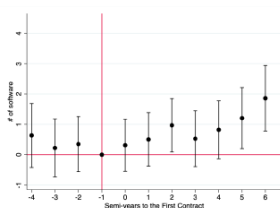
# Baseline empirical strategy

▶ **Triple diffs:** compare cumulative software releases before and after firms received 1st data-rich contracts, relative to the data-scarce ones

$$y_{it} = \sum_T \beta_{1T} T_{it} Data_i + \sum_T \beta_{2T} T_{it} + \alpha_t + \gamma_i + \sum_T \beta_{3T} T_{it} X_i + \epsilon_{it}$$

- $T_{it}$: 1 if, at time $t$, $T$ semi-years have passed before/since firm $i$ received 1st contract
- $Data_i$: 1 if firm $i$ receives "data rich" contract (i.e., from "high" camera capacity prefecture at time of contract receipt)
- $X_i$ controls for pre-contract firm characteristics: age, size (cap), and software production

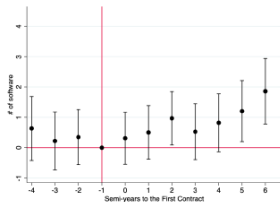# Public security contract "richer in data" & firm innovation

**Commercial use cumulative software releases**

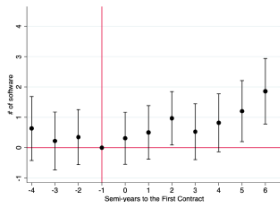# Public security contract "richer in data" & firm innovation

**Commercial use cumulative software releases**
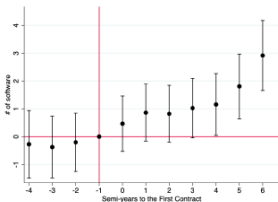


Magnitude: 2 new software products over 3 years
(20% of pre-contract software)

# Public security contract "richer in data" & firm innovation

**Commercial use cumulative software releases**



**Government use cumulative software releases**



Commercial innovation overcomes crowd-out of inputs by gov't

# Evaluating alternative hypotheses

1. **Selection** at a given time differs by contract
   - No differential pre-contract levels/trends of software
   - Control for time-varying effects of proxies for firms' underlying productivity: index constructed from establishment year, pre-contract capitalization, pre-contract rounds of external financing, pre-contract software production

# Evaluating alternative hypotheses

1. **Selection** at a given time differs by contract
   - No differential pre-contract levels/trends of software
   - Control for time-varying effects of proxies for firms' underlying productivity: index constructed from establishment year, pre-contract capitalization, pre-contract rounds of external financing, pre-contract software production

2. **Productive benefits other than data** differ by contract
   - Value of contract; tasks of contract; market access: we control for time-varying effects of an index of non-data contract characteristics (dollar value; prefecture income; tasks coded using NLP)
   - Signaling value: examine second contracts within parent firms
   - Political value: drop Beijing/Shanghai contracts; drop firms receiving contracts in home province
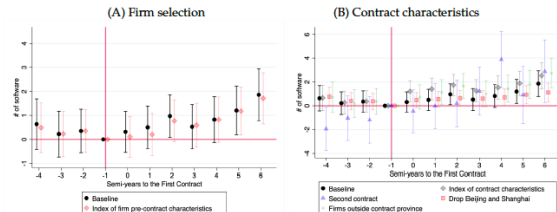
# Evaluating alternative hypotheses

1. **Selection** at a given time differs by contract
   - No differential pre-contract levels/trends of software
   - Control for time-varying effects of proxies for firms' underlying productivity: index constructed from establishment year, pre-contract capitalization, pre-contract rounds of external financing, pre-contract software production
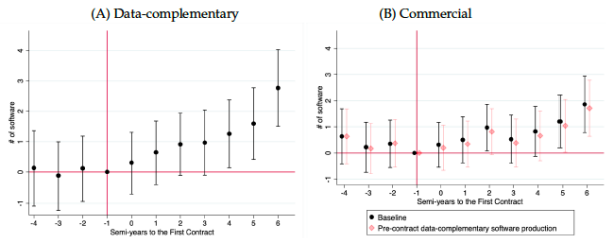
2. **Productive benefits other than data** differ by contract
   - Value of contract; tasks of contract; market access: we control for time-varying effects of an index of non-data contract characteristics (dollar value; prefecture income; tasks coded using NLP)
   - Signaling value: examine second contracts within parent firms
   - Political value: drop Beijing/Shanghai contracts; drop firms receiving contracts in home province

# Additional evidence for our mechanism(s)

**Data-complementary** software (e.g., storage/transmission) differentially increases after data-rich contract (learning); but, accounting for pre-contract data-complementary software does not greatly affect our findings (sharable data and algorithms)

# Contributions to literature

1. To the literature on the economics of AI and data (e.g., Aghion et al., 2017; Agrawal et al., 2018; Farboodi et al., 2019; Jones and Tonetti, 2019)

    - Highlight the role of **government data** in shaping commercial AI innovation, and the **sharability of data/algorithms** within the firm

# Contributions to literature

1. To the literature on the economics of AI and data (e.g., Aghion et al., 2017; Agrawal et al., 2018; Farboodi et al., 2019; Jones and Tonetti, 2019)

   - Highlight the role of **government data** in shaping commercial AI innovation, and the **sharability of data/algorithms** within the firm

2. To the literature on industrial and innovation policies (e.g., Rodrik, 2007; Lane, 2020; Bloom et al., 2019)

   - Government data provision to firms can act as an innovation policy, **whether intentionally or not**

   - Mechanisms **similar** to other government policies (e.g., learning spillovers from space exploration) but **distinct** too (direct effect of sharability)
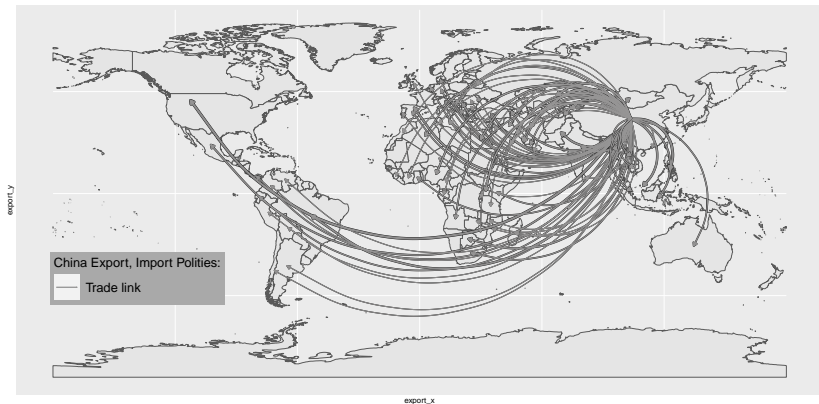
# Contributions to literature

1. To the literature on the economics of AI and data (e.g., Aghion et al., 2017; Agrawal et al., 2018; Farboodi et al., 2019; Jones and Tonetti, 2019)

   - Highlight the role of **government data** in shaping commercial AI innovation, and the **sharability of data/algorithms** within the firm

2. To the literature on industrial and innovation policies (e.g., Rodrik, 2007; Lane, 2020; Bloom et al., 2019)

   - Government data provision to firms can act as an innovation policy, **whether intentionally or not**

   - Mechanisms **similar** to other government policies (e.g., learning spillovers from space exploration) but **distinct** too (direct effect of sharability)

3. To the literature on the rise of China emphasizing the role of the state (e.g., Lau et al., 2000; Brandt and Rawski, 2008; Song et al., 2011)

   - Highlight the role of the **surveillance apparatus** in commercial innovation

   - *Next project:* **AI-tocracy**. Alignment between innovation and autocracy? Contrasts with e.g., North (1991); Acemoglu and Robinson (2006, 2012)
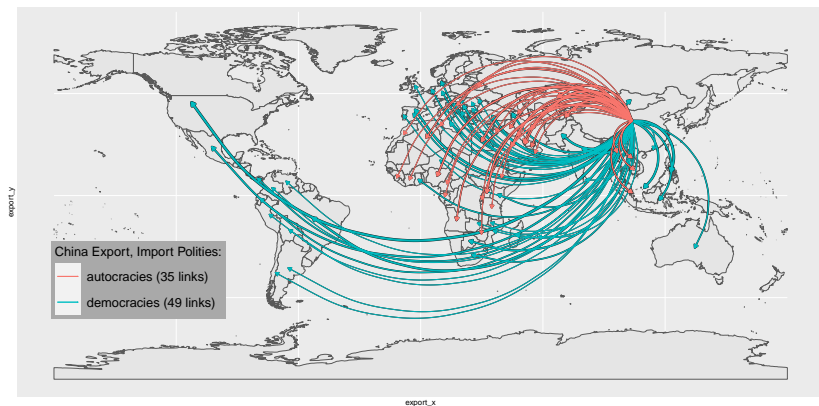
# China's export of AI

Dominate global trade (> 50%), different from other frontier tech

# China's export of AI

## High number of autocratic destinations

# US's export of AI

Much fewer links, higher share of democratic destinations