

# DATA, AI, AND THE STATE: EVIDENCE FROM CHINA

---

Martin Beraja (MIT)

NBER Digitization Tutorial, Spring 2022

# THE INTERPLAY BETWEEN DATA, AI, AND THE STATE

- ▶ Data-intensive technologies (like AI) can transform modern economies but have brought **new challenges** to the fore
- ▶ This has raised questions about **the role of governments**
  1. When technologies use data, what innovation policies and regulations are appropriate?
  2. Could governments misuse AI as a tool of repression and social control?

# THE INTERPLAY BETWEEN DATA, AI, AND THE STATE

- ▶ Data-intensive technologies (like AI) can transform modern economies but have brought **new challenges** to the fore
- ▶ This has raised questions about **the role of governments**
  1. When technologies use data, what innovation policies and regulations are appropriate?
  2. Could governments misuse AI as a tool of repression and social control?
- ▶ **Today:** think about the **interplay** between data, AI, and the state
  - ▶ Emphasize that AI is not only data-intensive but also **dual-use**
  - ▶ Prototypical setting to study this question: **the facial recognition AI industry in China**

1. Data-intensive Innovation and the State: Evidence from AI firms in China (with David Yang and Noam Yuchtman)
2. AI-tocracy (with Andrew Kao, David Yang and Noam Yuchtman)
3. Exporting Autocracy via AI (with Andrew Kao, David Yang and Noam Yuchtman)

## 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...

⇒ **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 this data can use it to train algorithms and develop software

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

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

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

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

⇒ 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

## 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)

## 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)

### 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)

## 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

- 10,677 AI contracts issued by public security arms of government (e.g., local police department)
- Data also on procurement of **AI-capable surveillance cameras**

中国政府采购网 首页 • 地方招采 • 中标公告

道路交通安全综合管理平台维护升级项目中标（成交）公告

2016年12月30日 16:28 来源：中国政府采购网 【打印】 【回到顶部】

1. 项目名称:道路交通安全综合管理平台维护升级项目  
2. 项目编号: GZGC-2016-38  
3. 项目序列号: S5200000000007081001  
4. 项目联系人: 王继刚  
5. 项目联系人电话: 0851-45226523  
6. 项目用途、简要技术要求及合同履行日期: 嵌入式“人脸识别”系统软件开发  
7. 采购方式: 公开招标  
8. 采购日期: 2016-12-07  
9. 公告媒体: 贵州省政府采购网  
10. 评审时间: 2016-12-29  
11. 评审地点: 贵州省公共资源交易中心  
12. 评审委员会成员名单:  
熊险峰、李强、田铁化、戚玉峰、袁荣伟  
13. 定标日期: 2016-12-29  
14. 中标（成交）信息:

序号	中标供应商	中标供应商地址	主要中标内容	中标金额 (元)
1	上海依图网络科技有限公司	上海市闵行区吴中路189号, 德必易创330-444室	嵌入式“人脸识别”系统软件开发	639000.00

15. 中标公告: 否

16. 采购人名称: 贵州省公安厅交通管理局  
联系地址: 贵阳市龙堡堡路416号  
项目联系人: 宋先生  
联系电话: 0851-45226880

17. 采购代理机构名称: 贵州贵财招标有限责任公司  
联系地址: 贵州省贵阳市观山湖区长岭北路233号贵州产业投资（集团）有限责任公司大楼413室  
项目联系人: 王继刚  
联系电话: 0851-45226523

18. 采购文件上传 (PDF格式):  
附件:  
gzc-2016-38(12月2日修改版).pdf

19. 书面推荐供应商参加采购活动的采购人和评审专家推荐意见 (如有):  
无

贵州贵财招标有限责任公司

**Deal Time** (points to item 14)

**Products/Services** (points to item 6)

**Monetary Scale** (points to item 14)

**Supplier** (points to item 16)

**Buyer** (points to item 16)



Registered with Min. of Industry and Information Technology

Categorize by intended customers (with RNN model using tensorflow):

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.*

**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

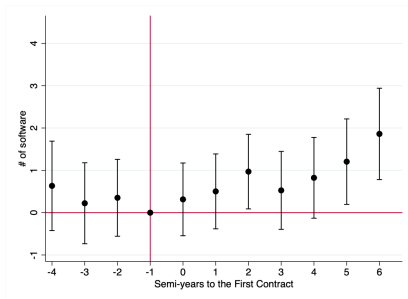
- **Triple diffs:** compare cumulative software releases before and after firms received 1st data-rich contracts, relative to the data-scarce ones

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

- $T_{it}$ : 1 if, at time  $t$ ,  $T$  semi-years have passed before/since firm  $i$  received 1st contract
- $\text{Data}_i$ : 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

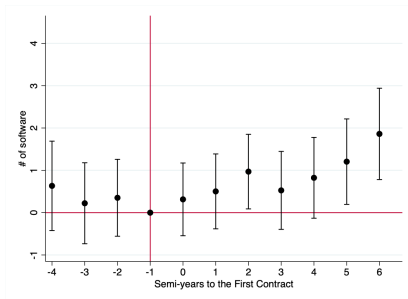
## Commercial cumulative software releases



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

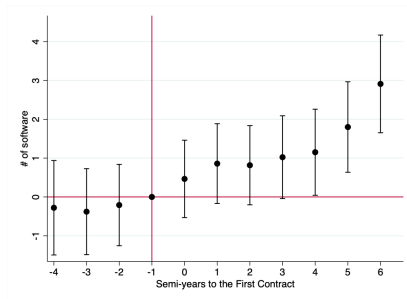
# PUBLIC SECURITY CONTRACT “RICHER IN DATA” & FIRM INNOVATION

## Commercial cumulative software releases



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

## Government cumulative software releases



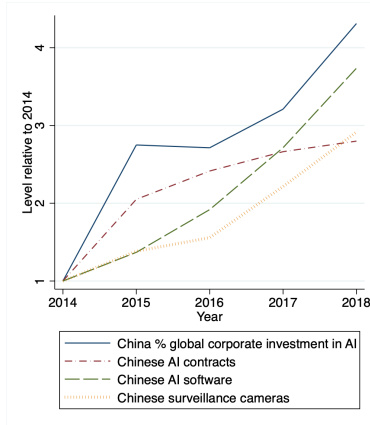
Commercial innov. overcomes gov't crowd-out

## IN THE PAPER: ASSESS MECHANISM(S) AND ALTERNATIVE HYPOTHESES

1. Selection at a given time differs by contract?
2. Productive benefits other than data differ by contract?
3. Data/algorithm sharability v. learning?

# TAKEAWAYS

1. Access to gov't data contributed to Chinese AI firms' emergence as leading innovators
  - Indeed, this has coincided with the expansion of the state's AI procurement and surveillance capacity



1. Access to gov't data contributed to Chinese AI firms' emergence as leading innovators
  - ▶ Indeed, this has coincided with the expansion of the state's AI procurement 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



1. Data-intensive Innovation and the State: Evidence from AI firms in China (with David Yang and Noam Yuchtman)
2. AI-tocracy (with Andrew Kao, David Yang and Noam Yuchtman)
3. Exporting Autocracy via AI (with Andrew Kao, David Yang and Noam Yuchtman)

## MOTIVATION: SUSTAINED INNOVATION UNDER AUTOCRACY?

## MOTIVATION: SUSTAINED INNOVATION UNDER AUTOCRACY?

- ▶ **Conventional wisdom:** autocracies are fundamentally misaligned with innovation  
(Lipset, 1959; Barro, 1996; Acemoglu and Robinson 2006; Glaeser et al., 2007; North et al. 2009)

## MOTIVATION: SUSTAINED INNOVATION UNDER AUTOCRACY?

- ▶ **Conventional wisdom:** autocracies are fundamentally misaligned with innovation (Lipset, 1959; Barro, 1996; Acemoglu and Robinson 2006; Glaeser et al., 2007; North et al. 2009)
- ▶ **This paper:** innovation in frontier techs. can be sustained under autocracy when they mutually reinforce each other

## MOTIVATION: SUSTAINED INNOVATION UNDER AUTOCRACY?

- ▶ **Conventional wisdom:** autocracies are fundamentally misaligned with innovation  
(Lipset, 1959; Barro, 1996; Acemoglu and Robinson 2006; Glaeser et al., 2007; North et al. 2009)
- ▶ **This paper:** innovation in frontier techs. can be sustained under autocracy when they mutually reinforce each other
  1. Frontier technology increases autocrats' probability of maintaining power
  2. Autocrats' spending on this tech. generates broader innovation spillovers

## MOTIVATION: SUSTAINED INNOVATION UNDER AUTOCRACY?

- ▶ **Conventional wisdom:** autocracies are fundamentally misaligned with innovation (Lipset, 1959; Barro, 1996; Acemoglu and Robinson 2006; Glaeser et al., 2007; North et al. 2009)
- ▶ **This paper:** innovation in frontier techs. can be sustained under autocracy when they mutually reinforce each other
  1. Frontier technology increases autocrats' probability of maintaining power
  2. Autocrats' spending on this tech. generates broader innovation spillovers
- ▶ AI may possess features that lead to a mutually reinforcing relationship

## MOTIVATION: SUSTAINED INNOVATION UNDER AUTOCRACY?

- ▶ **Conventional wisdom:** autocracies are fundamentally misaligned with innovation (Lipset, 1959; Barro, 1996; Acemoglu and Robinson 2006; Glaeser et al., 2007; North et al. 2009)
- ▶ **This paper:** innovation in frontier techs. can be sustained under autocracy when they mutually reinforce each other
  1. Frontier technology increases autocrats' probability of maintaining power
  2. Autocrats' spending on this tech. generates broader innovation spillovers
- ▶ **AI may possess features that lead to a mutually reinforcing relationship**
  1. As a technology of prediction, autocrats may be able to use AI for social / political control (Zuboff, 2019; Tirole, 2021; Acemoglu, 2021)
  2. Traditional spillovers (Moretti et al. 2019) + Sharability of gov't data/algorithms

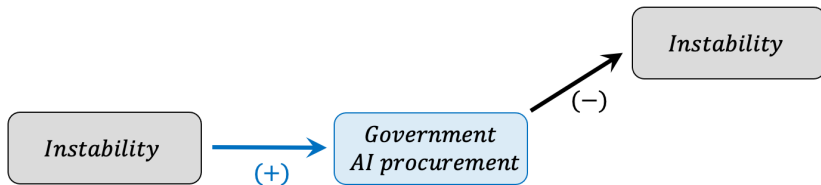
Test for a **mutually reinforcing relationship** between frontier innovation and autocracy in the context of China's facial recognition AI sector



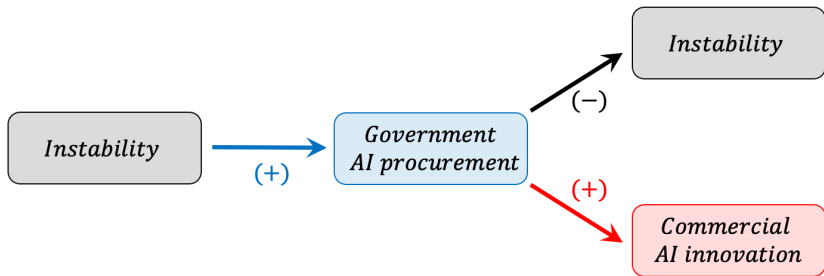
Test for a **mutually reinforcing relationship** between frontier innovation and autocracy in the context of China's facial recognition AI sector



Test for a **mutually reinforcing relationship** between frontier innovation and autocracy in the context of China's facial recognition AI sector



Test for a **mutually reinforcing relationship** between frontier innovation and autocracy in the context of China's facial recognition AI sector



### Protests and other episodes of political unrest:

- ▶ Daily level events in China from GDELT, a database tracking hundreds of news sites
- ▶ Use machine learning analysis to classify articles into those indicating political unrest (protests, demands, threats, etc.)
- ▶ There are 9,267 of these events from 2014 - 2020 throughout China

### Protests and other episodes of political unrest:

- ▶ Daily level events in China from GDELT, a database tracking hundreds of news sites
- ▶ Use machine learning analysis to classify articles into those indicating political unrest (protests, demands, threats, etc.)
- ▶ There are 9,267 of these events from 2014 - 2020 throughout China

### Weather

- ▶ Daily weather data from 260 weather stations across China
- ▶ LASSO regression to predict unrest events with 30 weather variables (e.g., temperature, precipitation, windspeed) and their interactions

# $\uparrow$ LOCAL UNREST IN QUARTER $t \implies \uparrow$ PUBLIC SECURITY AI PROCUREMENT IN $t + 1$

1. Diff-in-diff: panel specification, controlling for location and time FEs
2. IV: instrument unrest with [local weather conditions](#)
3. AI x Cameras: complementarity?

	<i>Public security AI procurement</i>			
	(1)	(2)	(3)	(4)
Panel A.1: OLS, AI				
Unrest events	0.199*** (0.043)	0.198*** (0.045)	0.199*** (0.044)	0.200*** (0.043)
Panel A.2: Lasso IV, AI				
Unrest events	0.388*** (0.088)	0.387*** (0.088)	0.388*** (0.088)	0.388*** (0.087)
Panel B.1: OLS, AI X surveillance cameras				
Unrest events	0.681*** (0.154)	0.669*** (0.157)	0.680*** (0.155)	0.674*** (0.150)
Panel B.2: Lasso IV, AI X surveillance cameras				
Unrest events	1.099*** (0.390)	1.083*** (0.385)	1.099*** (0.390)	1.085*** (0.384)
GDP $\times$ time	Yes	No	No	Yes
Population $\times$ time	No	Yes	No	Yes
Gov. revenue $\times$ time	No	No	Yes	Yes

## ↑ PUBLIC SECURITY AI STOCK IN QUARTER $t \implies$ EFFECT ON UNREST $t + 1$ ?

- Problematic to directly examine effect of AI stock on subsequent unrest events  
Positive autocorrelation between such events; and AI procurement is endogenous

# $\uparrow$ PUBLIC SECURITY AI STOCK IN QUARTER $t \implies \downarrow$ UNREST DUE TO GOOD WEATHER $t + 1$

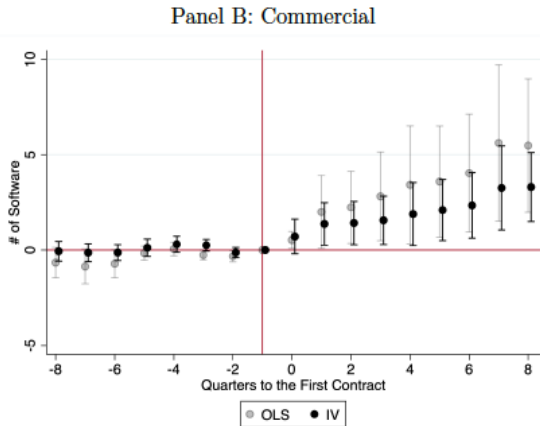
- Instead, examine whether AI tempers the effect of good weather on unrest events
- Also, look at AI in combinations with cameras, and placebo using non-public security AI

	<i>Standardized number of unrest events</i>							
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Panel A: Procurement of AI								
Favorable weather	0.9082*** (0.1576)	0.9422*** (0.1564)	0.9089*** (0.1579)	0.9410*** (0.1510)	0.9315*** (0.1646)	0.9705*** (0.1632)	0.9323*** (0.1650)	0.9684*** (0.1574)
Public security procurement stock $AI_{t-1}$	-0.0096** (0.0048)	-0.0057 (0.0061)	-0.0096** (0.0048)	-0.0044 (0.0056)				
Favorable weather $\times$ public security $AI_{t-1}$	-0.2626* (0.1563)	-0.3152* (0.1742)	-0.2623* (0.1570)	-0.3088* (0.1687)				
Non-public security procurement stock $AI_{t-1}$					-0.0025 (0.0017)	-0.0027 (0.0020)	-0.0025 (0.0017)	-0.0024 (0.0018)
Favorable weather $\times$ non-public security $AI_{t-1}$					-0.0492 (0.0367)	-0.0576 (0.0411)	-0.0495 (0.0372)	-0.0535 (0.0375)
Panel B: Procurement of AI X procurement of surveillance cameras								
Favorable weather	0.8989*** (0.1549)	0.9325*** (0.1524)	0.8994*** (0.1552)	0.9327*** (0.1480)	0.9554*** (0.1691)	0.9945*** (0.1659)	0.9562*** (0.1695)	0.9926*** (0.1605)
Public security procurement stock $AI_{t-1}$	0.2923*** (0.1083)	0.3158*** (0.0991)	0.2917*** (0.1081)	0.3081*** (0.0948)				
Favorable weather $\times$ public security $AI_{t-1}$	-0.7096*** (0.2302)	-0.7952*** (0.2412)	-0.7144*** (0.2323)	-0.7789*** (0.2248)				
Non-public security procurement stock $AI_{t-1}$					0.0605 (0.0600)	0.0626 (0.0592)	0.0608 (0.0603)	0.0601 (0.0572)
Favorable weather $\times$ non-public security $AI_{t-1}$					0.7558 (0.6020)	0.8049 (0.6015)	0.7573 (0.6043)	0.7744 (0.5801)
GDP $\times$ time	Yes	No	No	Yes	Yes	No	No	Yes
Log population $\times$ time	No	Yes	No	Yes	No	Yes	No	Yes
Gov. revenue $\times$ time	No	No	Yes	Yes	No	No	Yes	Yes



↑ POLITICALLY MOTIVATED PUBLIC SECURITY AI PROCUREMENT IN QUARTER  $t \implies$   
↑ COMMERCIAL AI INNOVATION IN  $t + 1$

1. Politically motivated public security contracts: those from location with above median unrest at  $t - 1$
2. Triple Diff: before/after firms receive 1st politically motivated contract, then compare to non-public sec. contracts



1. Alignment between autocrats demand for social control and AI innovation
  - ▶ Could shed light on prominent episodes of frontier innovation under non-democracies
    - ▶ Aerospace technology in the USSR
    - ▶ Chemical engineering innovation in Imperial Germany

## 1. Alignment between autocrats demand for social control and AI innovation

- ▶ Could shed light on prominent episodes of frontier innovation under non-democracies
  - ▶ Aerospace technology in the USSR
  - ▶ Chemical engineering innovation in Imperial Germany

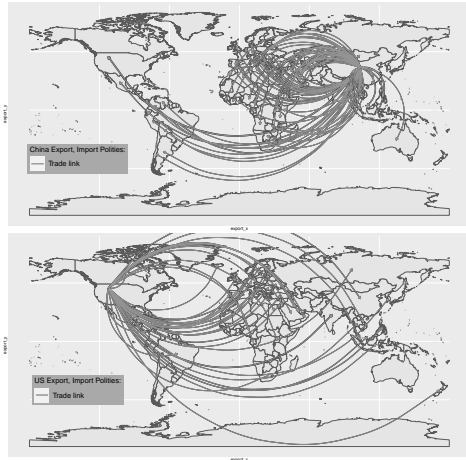
## 2. If China exports autocracy-enhancing AI, what are the international ramifications?

- ▶ China's comparative advantage in AI:  $\uparrow$  state demand  $\implies$   $\uparrow$  firms' global competitiveness
- ▶ Political bias: support autocracies and weak democracies abroad

1. Data-intensive Innovation and the State: Evidence from AI firms in China (with David Yang and Noam Yuchtman)
2. AI-tocracy (with Andrew Kao, David Yang and Noam Yuchtman)
3. Exporting Autocracy via AI (with Andrew Kao, David Yang and Noam Yuchtman)

# CHINA'S COMPARATIVE ADVANTAGE IN AI

Exports more AI than the US, particularly so when compared to other frontier technologies



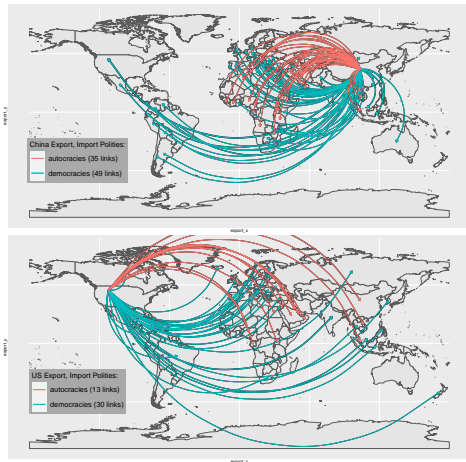
US vs. China, AI vs. frontier technologies

	<i>Linear probability of trade</i>			
	(1)	(2)	(3)	(4)
Origin China	-0.005 (0.004)	-0.005 (0.004)	0.005 (0.007)	-0.005 (0.004)
AI	-2.948*** (0.308)	-2.913*** (0.289)	-3.053*** (0.319)	-2.849*** (0.357)
Origin China X AI	0.220*** (0.045)	0.229*** (0.045)	0.262*** (0.061)	0.220*** (0.045)
N	5364	5364	5364	5364

*Notes:* Regressions are at the product-import-export country dyad level. Outcome is dummy for trade. Omitted: US X not AI. Errors clustered at origin countries. All columns control for import/export GDP and log distance. Column (2) adds controls for common border, free trade agreements, and shared colonial background. Column (3) adds controls for common language, legal system, and religion. Column (4) adds controls for landlocked and island characteristics.

# POLITICAL BIAS OF CHINESE AI EXPORTS

Imported more by autocracies and weak democracies. Not the case for the imports of China's other frontier technologies



AI vs. frontier technologies by polity type

	<i>Linear probability of trade</i>			
	(1)	(2)	(3)	(4)
Destination authoritarian	-0.043** (0.017)	-0.043** (0.017)	-0.046*** (0.016)	-0.041** (0.017)
Origin authoritarian	-0.107*** (0.026)	-0.107*** (0.026)	-0.110*** (0.025)	-0.105*** (0.025)
Origin China	-0.058*** (0.021)	-0.059*** (0.021)	-0.048** (0.021)	-0.054** (0.021)
AI	2.658*** (0.207)	2.657*** (0.209)	2.666*** (0.213)	2.668*** (0.212)
Destination authoritarian X AI	0.046*** (0.017)	0.046*** (0.017)	0.049*** (0.016)	0.044*** (0.017)
Origin authoritarian X AI	0.104*** (0.026)	0.104*** (0.026)	0.107*** (0.025)	0.102*** (0.026)
Origin China X AI	0.545*** (0.030)	0.547*** (0.033)	0.536*** (0.031)	0.542*** (0.031)
N	320796	320796	320796	320796

*Notes:* Regression at the product-import-export country dyad level. Outcome is dummy for trade. Omitted: origin/destination democracy X not AI. Errors two-way clustered at origin and destination countries. All columns control for import/export GDP and log distance. Column (2) adds controls for common border, free trade agreements, and shared colonial background. Column (3) adds controls for common language, legal system, and religion. Column (4) adds controls for landlocked and island characteristics.

# Thank you!

And if you want to chat more...  
maberaja@mit.edu