# Digitized data

DSIER [/dɪˈzaɪər/] — Summer 2023

Julian Hinz

Bielefeld University

# Today's plan

- "Non-computable information"
- Lloyd's shipping list: The Wind of Change: Maritime Technology, Trade, and Economic Development, Pascali (2017)
- Plantation records: "The Development Effects of the Extractive Colonial Economy: The Dutch Cultivation System in Java", Dell and Olken (2020)
- Clay tablets: "Trade, Merchants, and the Lost Cities of the Bronze Age", Barjamovic et al. (2019)

# NON-COMPUTABLE INFORMATION

# Non-computable information

- Standard digitization methods often fail to capture historical documents effectively

  $\rightarrow$ especially for less frequently used languages, scripts and settings

- Data may also be trapped in various types of images

- Text data contains a significant amount of non-computable information

# Economics and data

- Key economic questions necessitate disaggregated data: Misallocation, inequality, social mobility, welfare effects of trade

- Long-term digital disaggregated data uncommon

  $\rightarrow$ existing data predominantly originating from high resource contexts

- Growing academic interest, also due to much better computing power and methods

| Ground Truth Crop | EffOCR Localized Crop | Character Inner Product Similarity Rank | | | | |
|---|---|---|---|---|---|---|
| | | 1 | 2 | 3 | 4 | 5 |
| | | c | e | ( | C | L |
| | | A | n | R | : | { |
| | | o | v | c | e | l |
| | | f | r | t | { | Y |
| | | o | O | o | V | X |

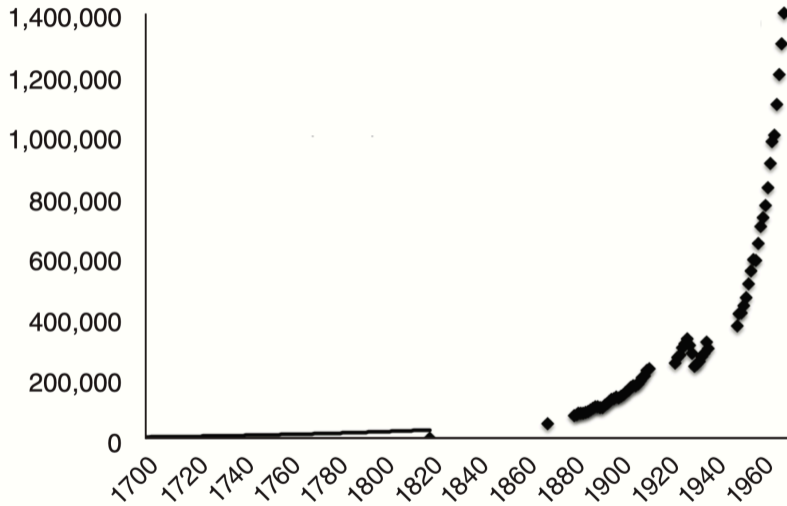| | | | | | | |
|---|---|---|---|---|---|---|
| 練 | 練 | 練 | **練** | 鍊 | 諫 | 涷 |
| 塚 | 塚 | 塚 | **塚** | 堠 | 彘 | 壖 |
| 麴 | 麴 | 魏 | **麴** | 麵 | 麲 | 麴 |
| 教 | 教 | 教 | 欻 | **教** | 資 | 譺 |
| | | 威 | 倣 | 倣 | **嫁** | 焱 |
| | | 豔 | **鹽** | 豔 | 鑑 | 鷸 |

# Accuracy

- OCR accuracy measured using character error rate (CER)

  → Levenshtein distance between recognized string and "ground truth", normalized by length of "ground truth"

  → minimum number of single-character edits (insertions, deletions or substitutions) required to change one word into the other

- CER of 0.5: mispredicting approximately half of characters
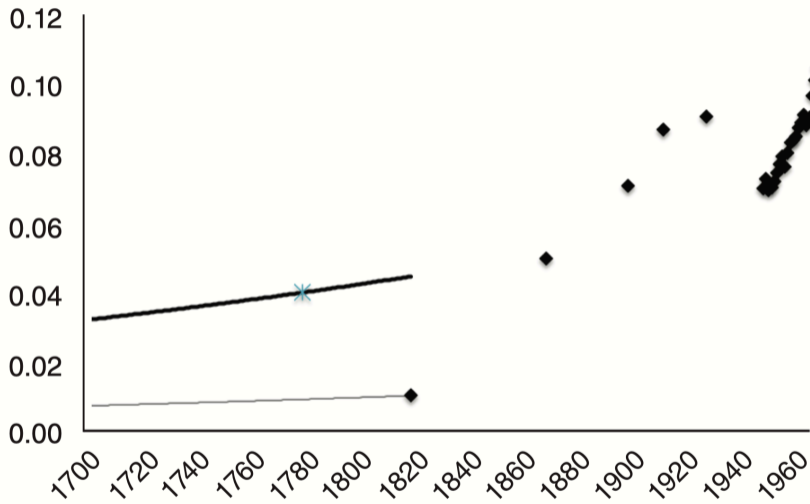
# Software solutions

- Google Cloud Vision, Amazon Textract

- Baidu OCR (for Asian languages)

- Tesseract (bi-directional LSTM)

- Active research: EasyOCR (Shi et al., 2016), TrOCR (Li et al., 2021), PaddleOCR (Du et al., 2022), ...
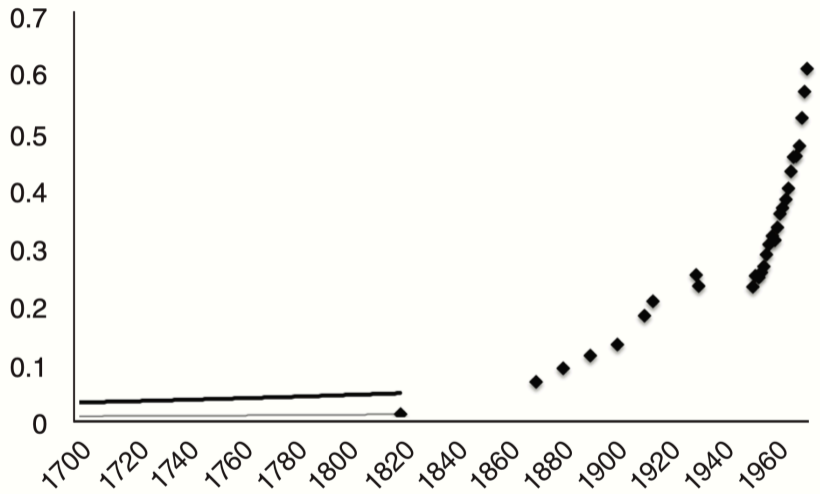
WIND OF CHANGE — LLOYD'S LIST

# Panel A. World exports at constant prices (million 1990 dollars)

Panel B. World export-to-GDP ratio

Panel C. World export-to-population ratio

# Idea

- 1870–1913 first era of trade globalization

- How did the increase in trade affect economic development?

- Causal mechanism: steamship vs. sailing

- asymmetric change in trade distances among countries

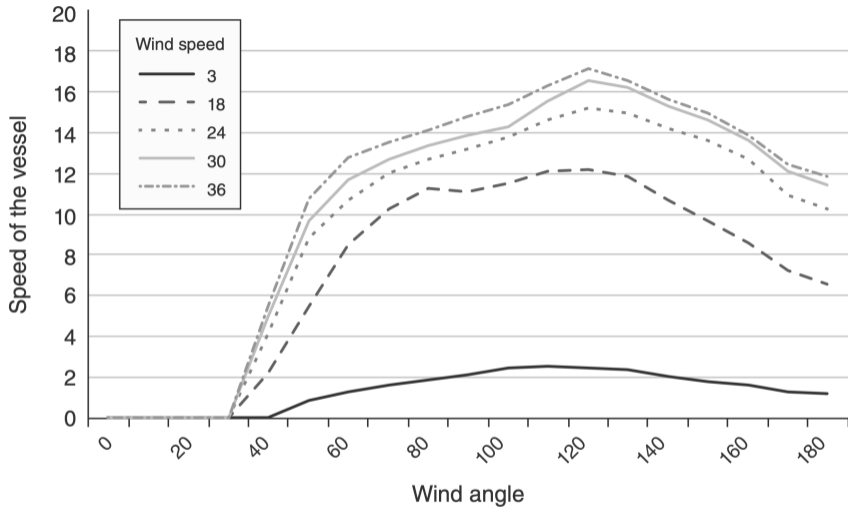- steamship reduced shipping costs and time heterogeneously across countries and trade routes
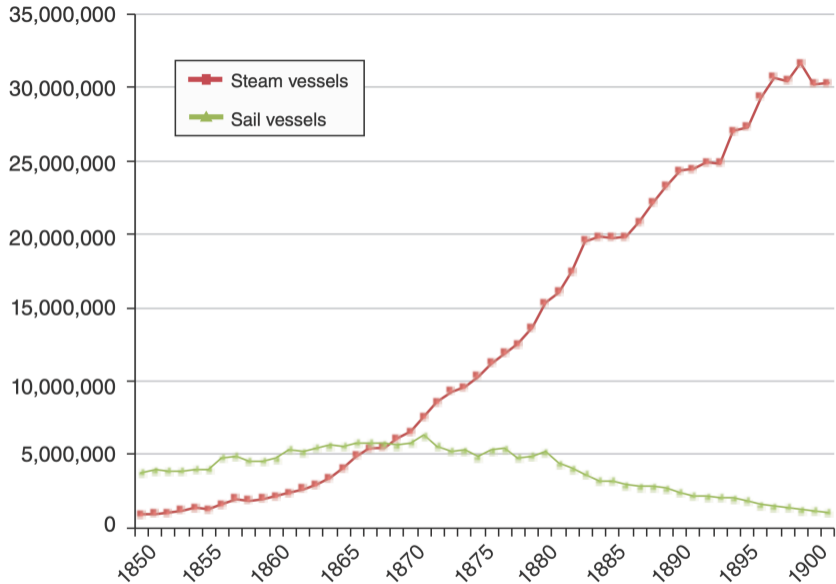
FIGURE 2. POLAR DIAGRAM OF A SAILING VESSEL: THE CLIPPER IN 1860
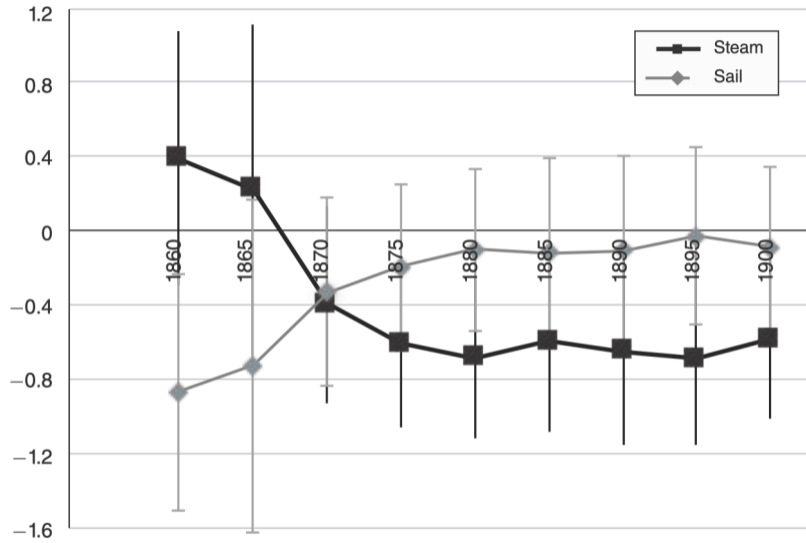
# Digitized data

- Three novel datasets from 1850 to 1900

- First dataset: shipping times across 16,000 country pairs

- Second dataset: 23,000 bilateral trade observations, 1,000 distinct country pairs

  $\rightarrow$ Sectoral-level export data for 37 countries

- Third dataset: freight rates across 291 shipping routes

# Effect on trade and GDP

- Impact of steamship on world trade volumes
  - $\rightarrow$ Reduction in geographical isolation measured by average shipping time
- Country-level regressions estimate impact of change in isolation

# Findings

- Rich countries did not benefit on average
- Similar impact of trade on agricultural and non-agricultural countries
- Institutions might reflect economic development differences

red

- Investigates Dutch Cultivation System impacts

- Farmers forced to cultivate export crops: Sugar

- Areas near factories more industrialized today

- Residents near factories have higher education

**Vervolg C.** Dessa's aan de onderneming dienstbaar, voor arbeid.

| Namen der Dessa's | Afstand in palen van de | | Kultuur-dienstplichtige huisgezinnen. | Beamankariet door ieder dessa te onderhouden. | Geeft koered huisgezinnen per boum. | Verwijzing naar de bedrieling(en). |
|---|---|---|---|---|---|---|
| | suikerriet-velden. | fabriek. | | | | |
| Transport | | | 402 | 63¾ | | |
| | | | | | | |

| Namen der Dessa's | Afstand in palen van de | | Kultuur-dienstplichtige huisgezinnen. | Beamankariet door ieder dessa te onderhouden. | Geeft koered huisgezinnen per boum. | Verwijzing naar de bedrieling(en). |
|---|---|---|---|---|---|---|
| | suikerriet-velden. | fabriek. | | | | |
| Transport | | | 983 | 148½ | | |
| Transporteren | | | 1996 | 224¾ | | |

- Data combines historical and contemporary sources

- Traces long-term impacts of Cultivation System

- Geographic distance to factories measures exposure

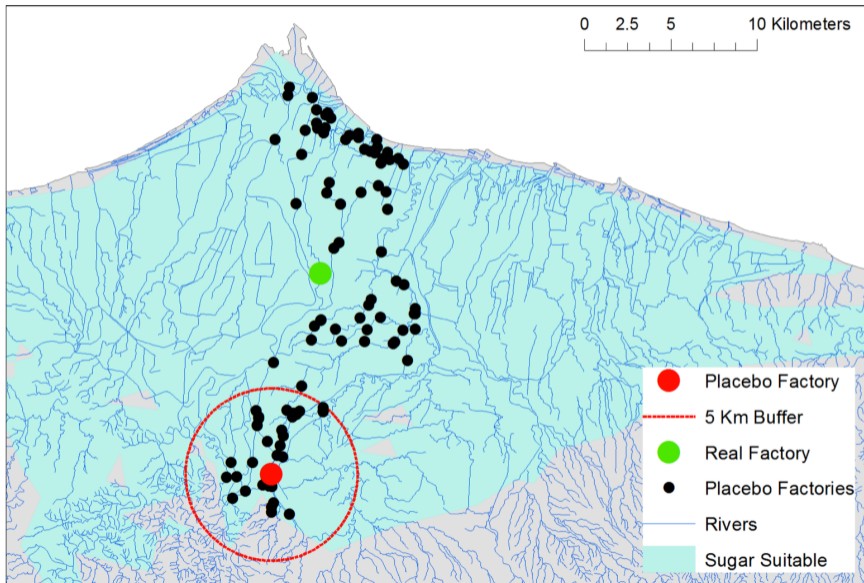- Uses contemporary data for long-term impacts
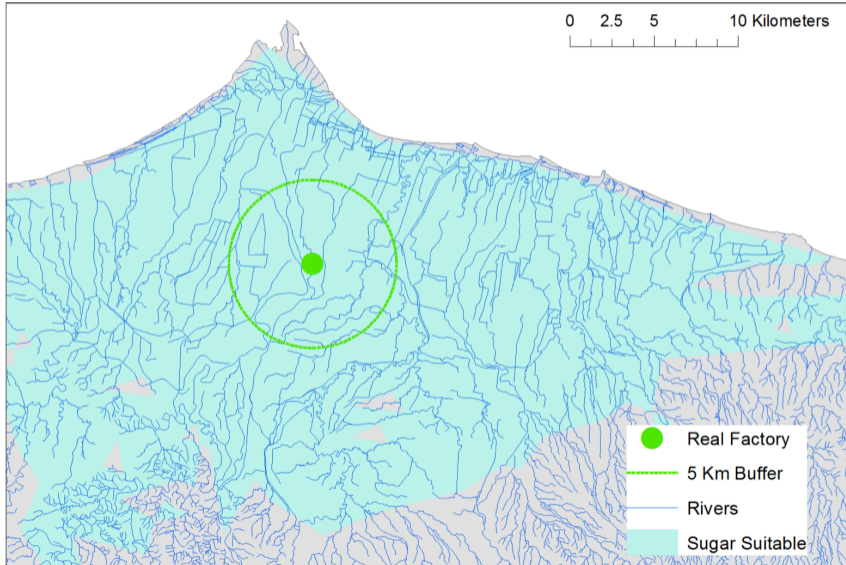
# (a) Real Factory

# (b) Placebo Factory Suitability

(c) Placebo Factories

# (a) Real Factory

# Discussion

- Study focuses on specific colonial institution

- Findings may not generalize to other institutions

- Pre-existing differences between areas not ruled out

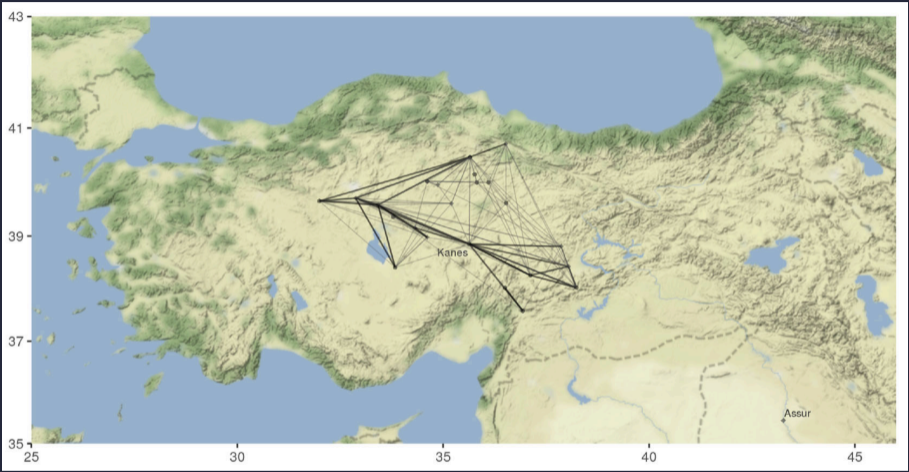- Unobserved factors could influence results

red

- Novel approach to estimate the locations of lost cities from the Bronze Age

- Structural gravity model to estimate the locations of lost cities based on trade data from ancient texts

- Ancient city sizes are persistent, meaning that large ancient cities tend to be located at or near large modern cities

# Data and its Novelty

- Sample of 9,728 digitized texts and approximately 2,000 additional non-digitized texts

- ancient texts to extract information about trade routes and city locations

- data mentions 79 unique settlements, with the analysis restricted to 25 Anatolian cities in Turkey

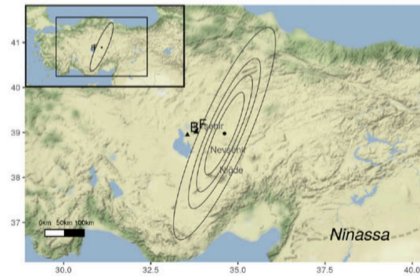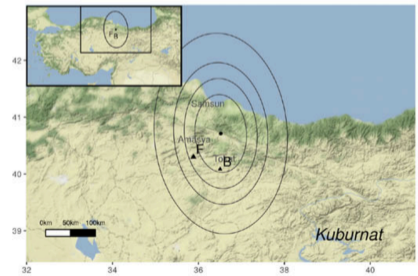FIGURE I
Tablet Kt 83-k 117

# Empirical Strategies

- structural gravity model to estimate the locations of lost cities

- detailed data on the topography of the entire region surrounding Anatolia to compute travel times

$$(4) \quad \Pr\left[c_{ij}(\omega) \leqslant \min_{k \in \mathcal{K}\backslash\{j\}}\left\{c_{kj}(\omega)\right\} \, \middle| \, \min_{l \in \mathcal{L} \cup \{j\}} c_{lj}(\omega) > \min_{k \in \mathcal{K}\backslash\{j\}}\left\{c_{kj}(\omega)\right\}\right]$$

$$= \frac{T_i \left(\tau_{ij} w_i\right)^{-\theta}}{\sum_{k \in \mathcal{K}\backslash\{j\}} T_k \left(\tau_{kj} w_k\right)^{-\theta}},$$

## TABLE II
### Lost Cities' Geocoordinates

|  | Latitude | (Std. err.) | Longitude | (Std. err.) | Correlation |
|---|---|---|---|---|---|
| Durhumit | 40.47 | (0.025) | 35.65 | (0.445) | − 0.952 |
| Hahhum | 38.429 | (0.274) | 38.04 | (0.517) | 0.68 |
| Kuburnat | 40.712 | (0.582) | 36.52 | (0.512) | − 0.06 |
| Ninašša | 38.977 | (0.778) | 34.614 | (0.482) | 0.86 |
| Purušhaddum | 39.71 | (1.54) | 32.872 | (0.669) | 0.774 |
| Šinahuttum | 39.956 | (0.333) | 34.866 | (0.165) | 0.863 |
| Šuppiluliya | 40.021 | (1,022.82) | 34.618 | (58.796) | 1.0 |
| Tuhpiya | 39.611 | (0.18) | 35.199 | (0.307) | 0.528 |
| Wašhaniya | 39.157 | (0.219) | 34.311 | (0.265) | − 0.01 |
| Zalpa | 38.805 | (0.648) | 37.862 | (1.199) | 0.878 |

*Notes*. This table presents the estimated geocoordinates, latitudes and longitudes, from solving our structural gravity model (8). All latitudes are north, and all longitudes are east. Robust (White) standard errors are in parentheses. The last column gives the estimated correlation between latitude and longitude, used to compute confidence regions.

# Discussion

- may be a systematic bias for larger cities to be more or less likely to have been unambiguously located by historians

- large ancient cities may never be discovered, as they lay buried under modern cities

- data does not observe internal transactions, a purchase in a city of a good sourced locally in the same city