



MENS  
MANUS AND  
MACHINA

# Tensorial Machine Learning for Discovering Patterns and Causality from Spatiotemporal Data

International Merchandise Trade & Urban Human Mobility

**Xinyu Chen**

M3S Project

May 20, 2024

# Outline

---

## A quick look:

- Motivations
- Frameworks
  - Pattern discovery
  - Causal effect imputation
- Spatiotemporal data
  - International merchandise trade
  - Urban human mobility
- Application to M3S

# Motivation

---

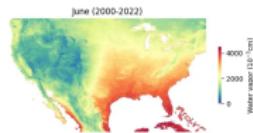
- Spatiotemporal systems & data scenarios



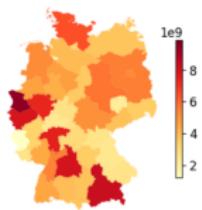
Transportation



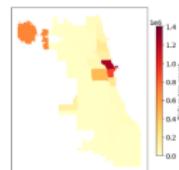
International trade



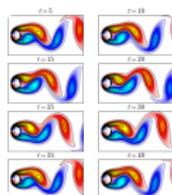
Climate



Energy



Mobility

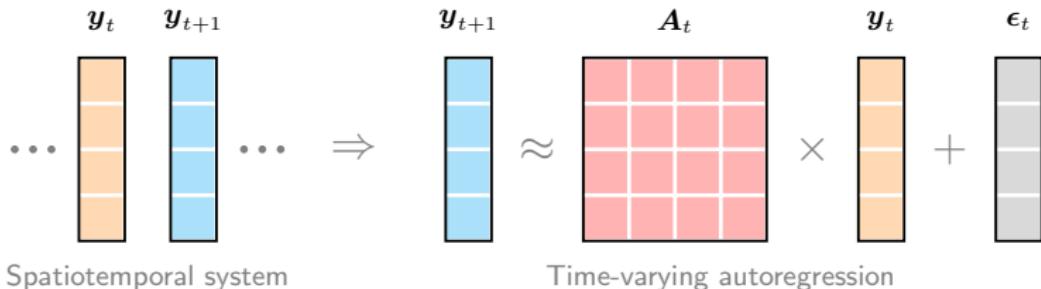


Fluid flow

# Pattern Discovery

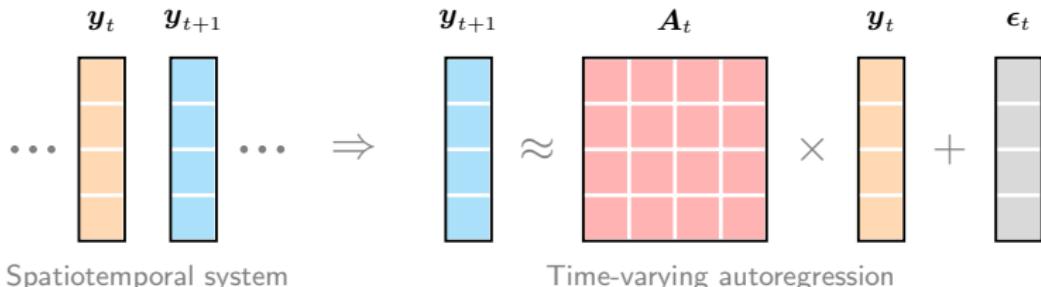
---

- How to characterize dynamical systems?



# Pattern Discovery

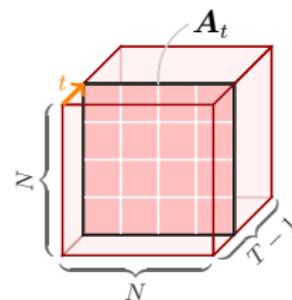
- How to characterize dynamical systems?

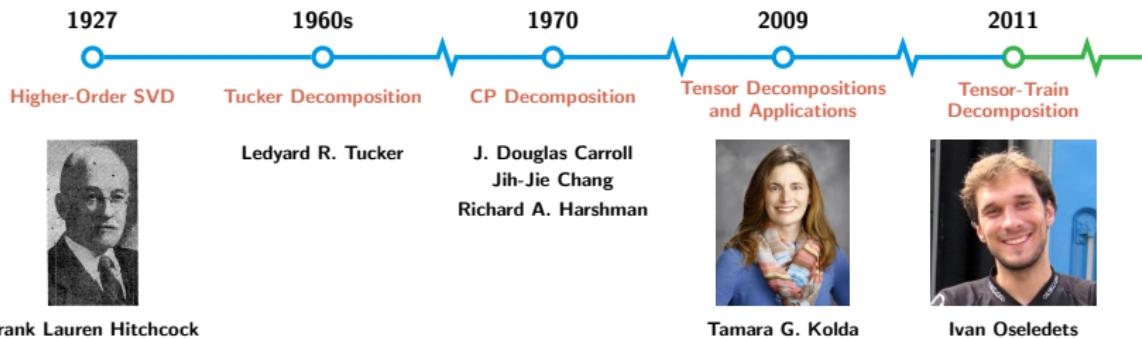


- On spatiotemporal systems  $\mathbf{Y} \in \mathbb{R}^{N \times T}$ :

$$\underbrace{y_{t+1} = Ay_t + \epsilon_t}_{\text{time-invariant (e.g., DMD)}} \quad \text{v.s.} \quad \underbrace{y_{t+1} = A_t y_t + \epsilon_t}_{\text{fully time-varying (ours)}}$$

- How to discover spatial/temporal modes (patterns) from the tensor  $\mathcal{A} \triangleq \{A_t\}_{t \in [T-1]}$ ?

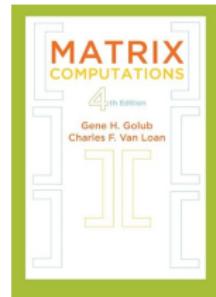




# Pattern Discovery

- Tensor factorization:

$$\mathcal{A} = \underbrace{\mathcal{G} \times_1 \mathbf{W} \times_2 \mathbf{V} \times_3 \mathbf{X}}_{\text{Tucker decomposition}}$$
$$\Downarrow$$
$$\mathbf{A}_t = \mathcal{G} \times_1 \underbrace{\mathbf{W}}_{\text{spatial modes}} \times_2 \mathbf{V} \times_3 \underbrace{\mathbf{x}_t^\top}_{\text{temporal modes}}$$



- (Ours)** Dynamic autoregressive tensor factorization (DATF):

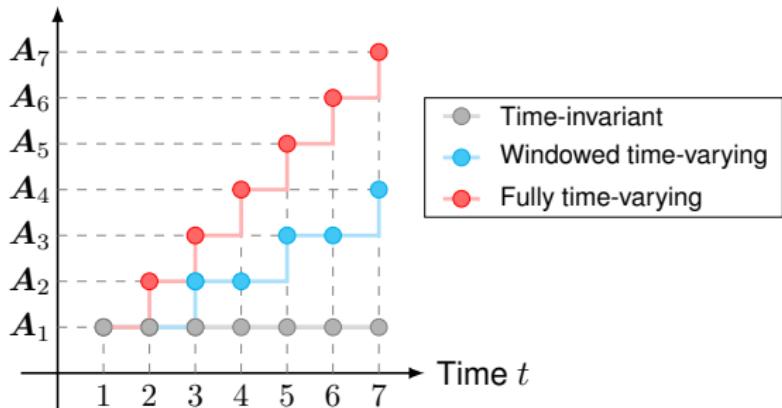
$$\min_{\mathcal{G}, \mathbf{W}, \mathbf{V}, \mathbf{x}} \frac{1}{2} \sum_{t \in [T-1]} \|\mathbf{y}_{t+1} - (\mathcal{G} \times_1 \mathbf{W} \times_2 \mathbf{V} \times_3 \mathbf{x}_t^\top) \mathbf{y}_t\|_2^2$$

s.t.  $\underbrace{\mathbf{W}^\top \mathbf{W} = \mathbf{I}_R}_{\text{orthogonal spatial modes}}$

- On the data  $\mathbf{Y} \in \mathbb{R}^{N \times T}$ :

$$\underbrace{\mathbf{y}_{t+1} = \mathbf{A}\mathbf{y}_t + \epsilon_t}_{\text{time-invariant (e.g., DMD)}} \quad \text{v.s.} \quad \underbrace{\mathbf{y}_{t+1} = \mathbf{A}_t \mathbf{y}_t + \epsilon_t}_{\text{fully time-varying (ours)}}$$

Coefficients



# Pattern Discovery

- Import/Export merchandise trade values (annual)<sup>1</sup> (215 countries/regions & period of 2000–2022)
  - Total merchandise trade values
  - Represent import/export trade data as a 215-by-23 matrix



Imports from 2000 to 2022



Exports from 2000 to 2022

<sup>1</sup>The dataset is available at <https://stats.wto.org>.



Import pattern 1



Export pattern 1



Import pattern 2



Export pattern 2



Import pattern 3



Export pattern 3



Import pattern 4

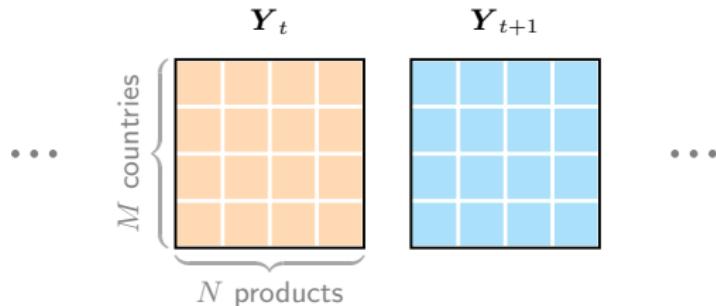


Export pattern 4

# Pattern Discovery

---

- Three-dimensional trade (Country, Product, Year)

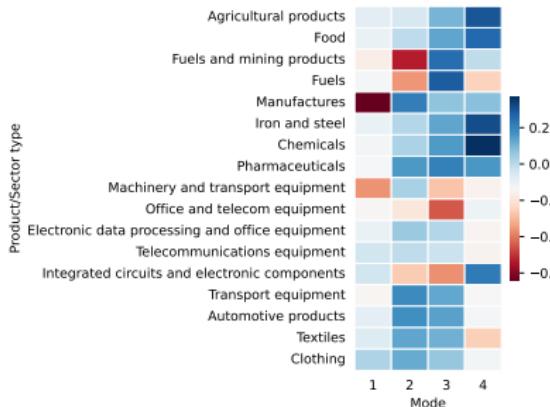


- On spatiotemporal systems  $\mathcal{Y} \in \mathbb{R}^{M \times N \times T}$ :

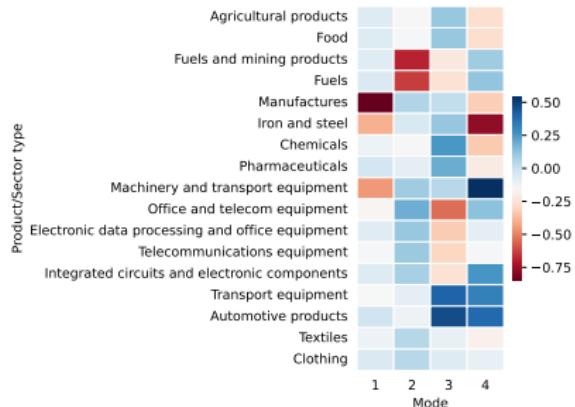
$$\underbrace{\mathbf{y}_{n,t+1} = \mathbf{A}_{n,t} \mathbf{y}_{n,t} + \boldsymbol{\epsilon}_{n,t}}_{\text{time-varying \& product-varying}}$$

# Pattern Discovery

- Product patterns on 17 merchandise types



Imports



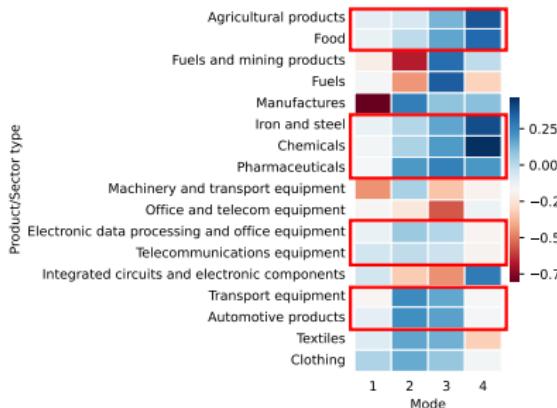
Exports

- Classify import/export merchandise according to product patterns
- Basic principle:

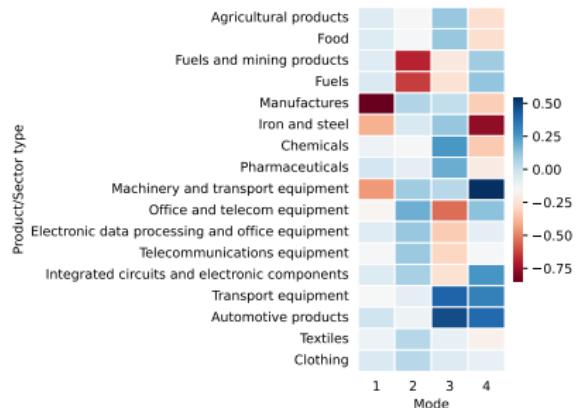
Import: What we buy? (demand) vs. Export: What we sell? (supply)

# Pattern Discovery

- Product patterns on 17 merchandise types



Imports



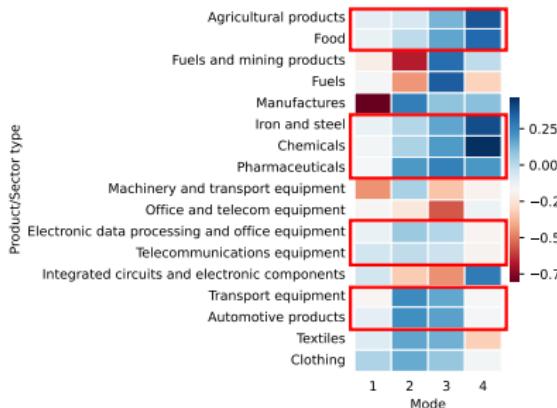
Exports

- Classify import/export merchandise according to product patterns
- Basic principle:

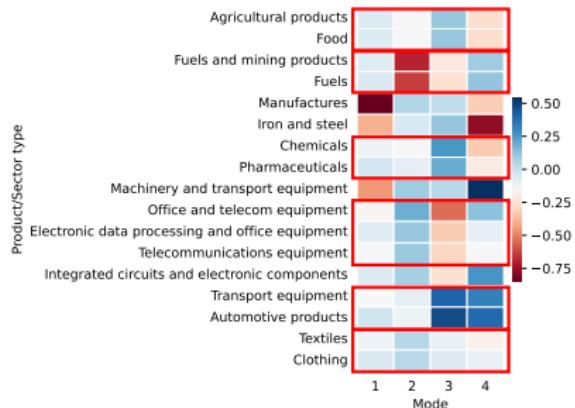
Import: What we buy? (demand) vs. Export: What we sell? (supply)

# Pattern Discovery

- Product patterns on 17 merchandise types



Imports



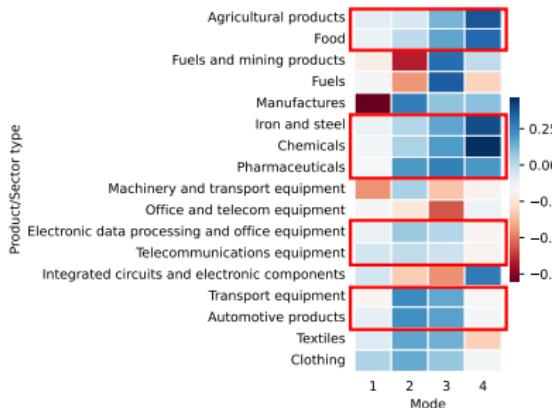
Exports

- Classify import/export merchandise according to product patterns
- Basic principle:

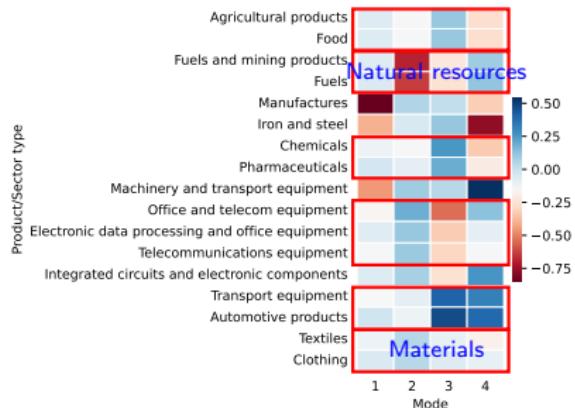
Import: What we buy? (demand) vs. Export: What we sell? (supply)

# Pattern Discovery

- Product patterns on 17 merchandise types



Imports



Exports

- Classify import/export merchandise according to product patterns
- Basic principle:

Import: What we buy? (demand) vs. Export: What we sell? (supply)

# Big International Trade Dataset

---

- Publicly available on Dropbox<sup>2</sup>
- Complicated dimensions<sup>3</sup>

Dimensions	Attributes
Exporter ID	231
Importer ID	231
Harmonized System (HS) code	5,018
Year	1995–2022 (28 years)
Value	Dollar (\$)
Quantity	Metric Tons (MT)

Note: 247,647,741 rows in total.

- Large data size: **3.5GB** zipped vs. **28GB** unzipped
- A lot of challenges unsolved:
  - Complicated dimensions, system evolution, non-negativity, long tails, high-order graphs/networks, causal effects, etc.

---

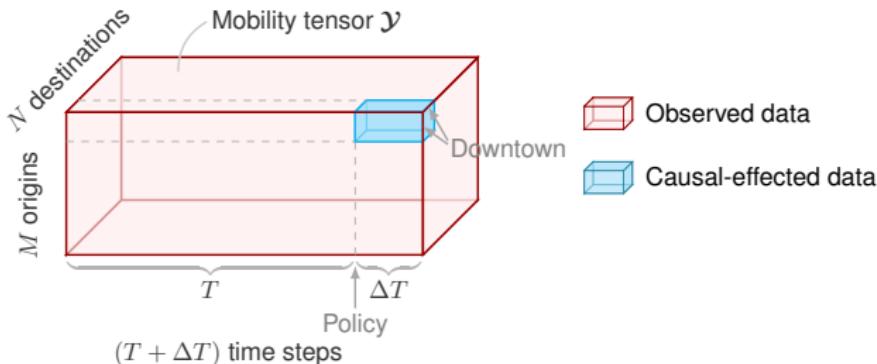
<sup>2</sup>Download the dataset at

[https://www.dropbox.com/s/eodmonwvvo25jkm/trade\\_i\\_baci\\_a\\_92.tsv.bz2?dl=0](https://www.dropbox.com/s/eodmonwvvo25jkm/trade_i_baci_a_92.tsv.bz2?dl=0)

<sup>3</sup>The HS is a standardized numerical method of classifying traded products.

# Causal Effect Imputation

- Urban human mobility w/ ODT
- Causal effect imputation framework

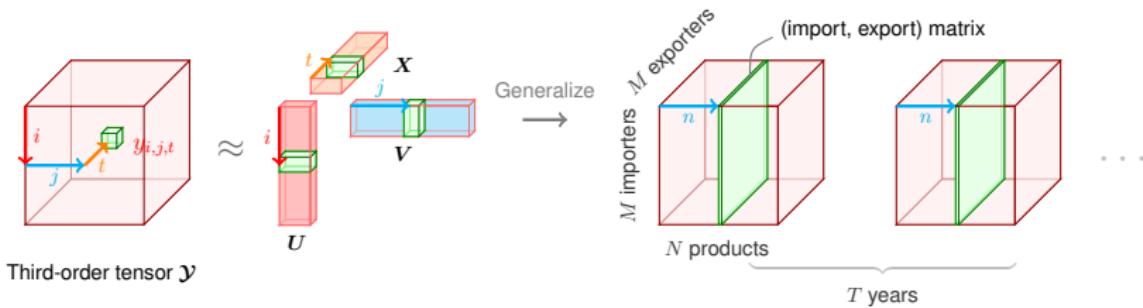


(Figure) Policy intervention (e.g., congestion pricing) in Downtown areas.

- A lot of challenges unsolved:
  - Structural missing patterns & cross validation
  - How to reduce data biases (e.g., underlying counterfactual factors) and model biases (e.g., utilizing spatial correlations & temporal dynamics)?
  - How to measure the difference between imputation results and realistic data?

# Causal Effect Imputation

- Not only for **mobility**, but also for **international trade**



- Technical challenges:
  - Nonconvex optimization
  - How to integrate spatiotemporal context (e.g., graph learning)?
  - Non-negativity ...

## Human Mobility Dataset

- Chicago taxi/ridesharing data

## Matching Taxi Trips with Community Areas

There are three basic steps to follow for processing test trip data:

- Download taxi trips in 2022 in the .csv format, e.g., `taxi_trips_-_2022.csv`
  - Use the `pandas` package in Python to process the raw trip data.
  - Match trip origin and destination locations with boundaries of the community area

```
import pandas as pd

data = pd.read_csv('Shel_Steps_-_2022.csv')
data.head()
```

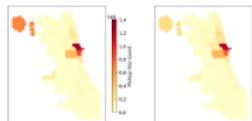
For each taxi trip, one can select some important information:

- **Trip Start Timestamp:** When the trip started, rounded to the nearest 15 minutes.
  - **Trip Seconds:** Time of the trip in seconds.
  - **Trip Miles:** Distance of the trip in miles.
  - **Pickup Community Area:** The Community Area where the trip began. This column will be blank for locations outside Chicago.
  - **Dropoff Community Area:** The Community Area where the trip ended. This column will be blank for locations outside Chicago.

```
df = pd.DataFrame()
df['Trip Start Timestamp'] = data['Trip Start Timestamp']
df['Trip Seconds'] = data['Trip Seconds']
df['Trip Miles'] = data['Trip Miles']
df['Pickup Community Area'] = data['Pickup Community Area']
df['Dropoff Community Area'] = data['Dropoff Community Area']

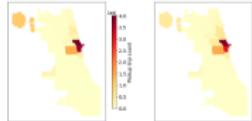
del data
del df
```

Figure 2 shows taxi pickup and dropoff trips (2022) on 77 community areas in the City of Chicago. Note that the average trip duration is 1207.75 seconds and the average trip distance is 8.16 miles.

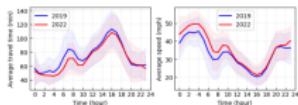


**Figure 2.** Taxi pickup and dropoff trips (2022) in the City of Chicago, USA. There are 4,763,961 remaining trips after the data processing.

For comparison, Figure 3 shows taxi pickup and dropoff trips (2019) on 77 community areas in the City of Chicago. Note that the average trip duration is 915.62 seconds and the average trip distance is 3.93 miles.



**Figure 3.** Taxi pickup and dropoff trips (2019) in the City of Chicago, USA. There are 12,484,572 remaining trips after the data processing. See the data processing code.



**Figure 6.** Average travel time and speed from area 32 (i.e., Downtown) to area 76 (i.e., Airport) in both 2019 and 2022.

Source: <https://spatiotemporal-data.github.io/Chicago-mobility/taxi-data>

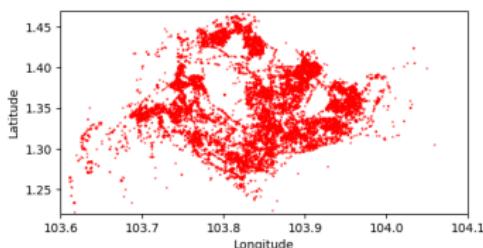
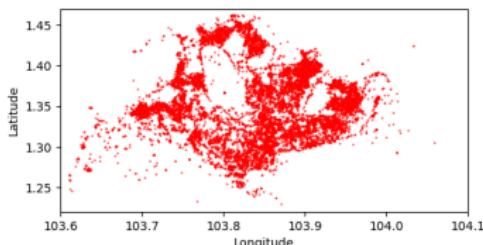
- Causal effect estimation of congestion pricing in Chicago

# Application to M3S

---

Veraset dataset in Singapore<sup>4</sup>

- Urban trajectory (samples)



- Trajectory inference (e.g., trip purpose identification)
- Pattern discovery
  - (On 2D activity data)  
Uncover spatial patterns  
(e.g., POI patterns)
  - (On 3D mobility data)  
Uncover temporal patterns  
(e.g., long-term pattern transition impacted by special events/policies)
- Causal effect imputation
  - Infer the causal effects of policy intervention

---

<sup>4</sup><https://spatiotemporal-data.github.io/trajectory/veraset/>



MENS  
MANUS AND  
MACHINA

# Thanks for your attention!

Any Questions?

## About me:

- 🏡 Homepage: <https://xinychen.github.io>
- ✉ How to reach me: [chenxy346@gmail.com](mailto:chenxy346@gmail.com)
- ✉ Or send to: [xinychen@mit.edu](mailto:xinychen@mit.edu)