# Data Science for Organizational Modeling
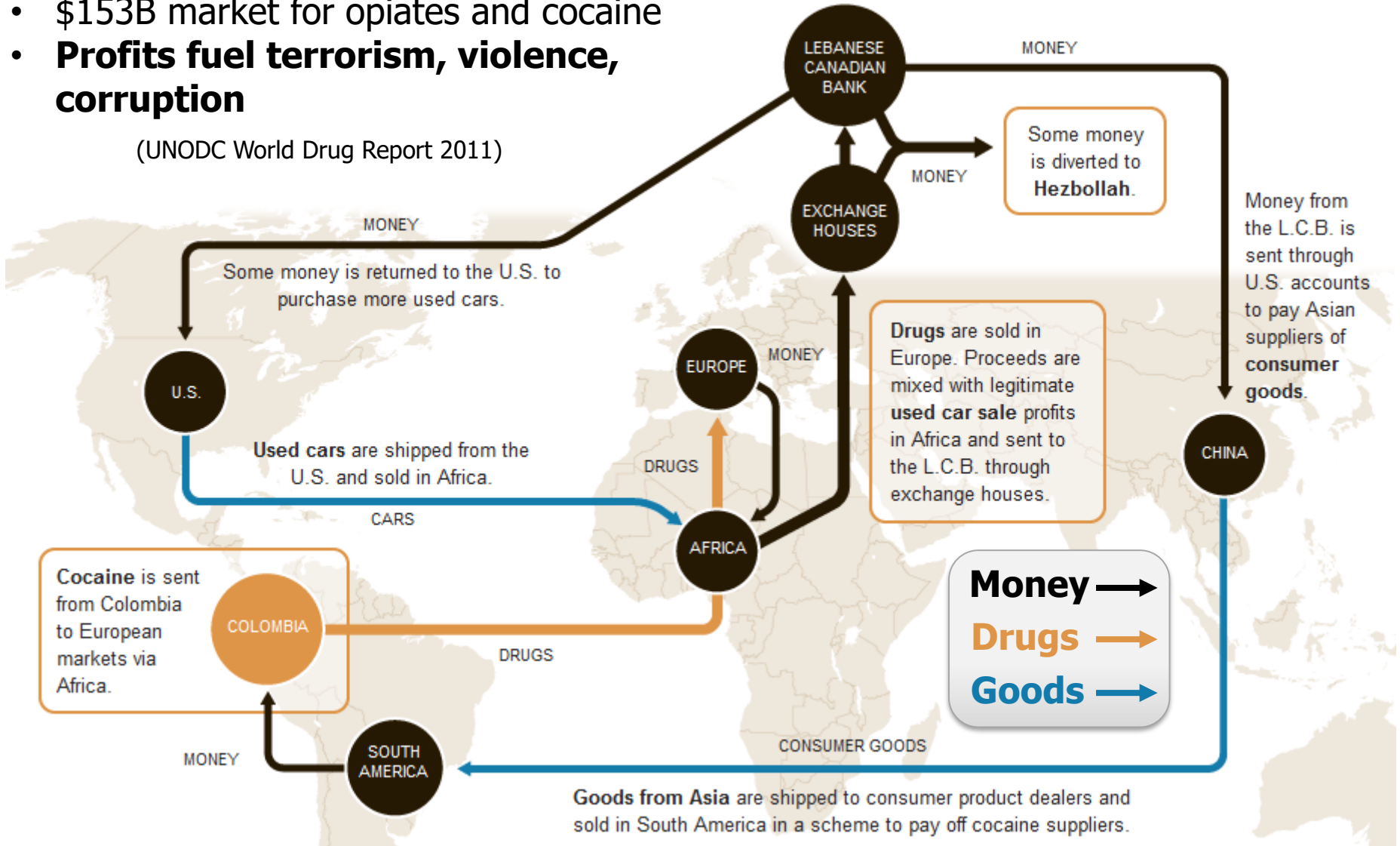
# Understanding Organizations and Their Relationships



**Traditional DOD (Nation State)**

**Unexplored Territory**

**Historic Military Counter-Insurgency**

- $153B market for opiates and cocaine
- **Profits fuel terrorism, violence, corruption**
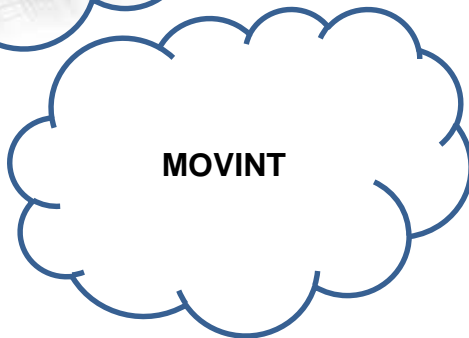
  (UNODC World Drug Report 2011)



LEBANESE CANADIAN BANK

MONEY

Some money is diverted to **Hezbollah.**

Money from the L.C.B. is sent through U.S. accounts to pay Asian suppliers of **consumer goods.**

EXCHANGE HOUSES

MONEY

MONEY

Some money is returned to the U.S. to purchase more used cars.

**Drugs** are sold in Europe. Proceeds are mixed with legitimate **used car sale** profits in Africa and sent to the L.C.B. through exchange houses.

EUROPE

MONEY

U.S.

**Used cars** are shipped from the U.S. and sold in Africa.

CARS

DRUGS

CHINA

AFRICA

**Cocaine** is sent from Colombia to European markets via Africa.

COLOMBIA

DRUGS

**Money** �ντ

**Drugs** ➡

**Goods** ➡

MONEY

SOUTH AMERICA

CONSUMER GOODS

**Goods from Asia** are shipped to consumer product dealers and sold in South America in a scheme to pay off cocaine suppliers.

Source: New York Times 2011    3

1. **Traditional Intel**
2. **Law enforcement**
3. **Online data**
4. **Commercial**

**Social Media**

**HUMINT**

**Cyber**

**Criminal Histories**

Overlapping data creates hyper-local observations and strategic insights not available through one source or one resolution alone

**Public Records**
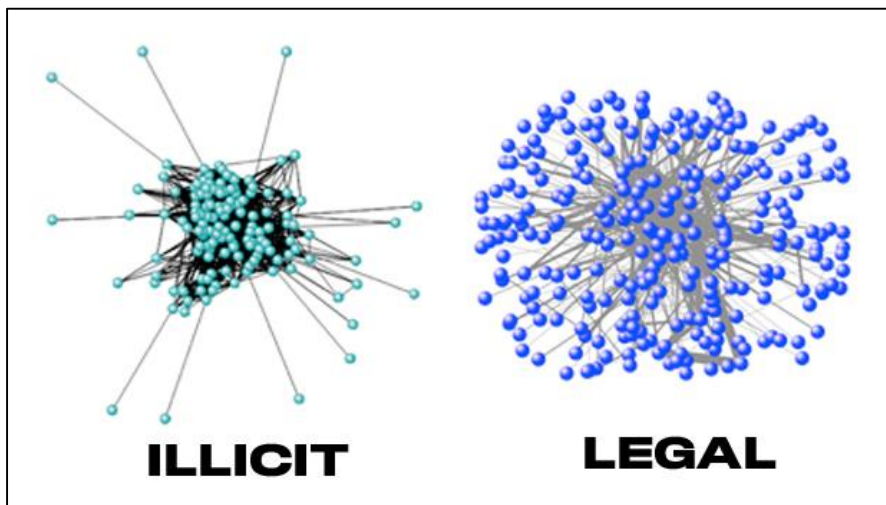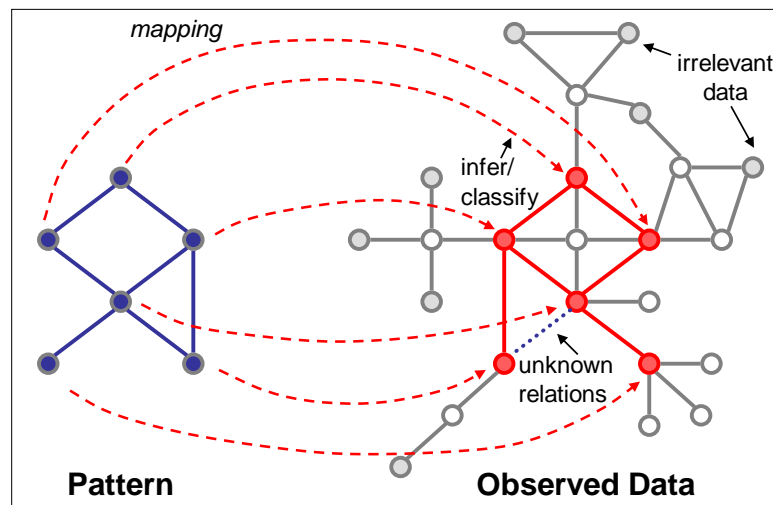
**Comms**

**News Feeds**

**MOVINT**

**FININT**

Detect exact and inexact patterns in networks to determine when predicted organizational behaviors match observations.

$$\min P(\mathbf{S}|\mathbf{D},\mathbf{M}) \cong \frac{1}{Z(\mathbf{D},\mathbf{M})} \prod_{ki}\left[p(a_{i,i}^{\mathbf{D}}|a_{k,k}^{\mathbf{M}})\right]^{S_{ki}} \prod_{kmij}\left[p(a_{i,j}^{\mathbf{D}}|a_{k,m}^{\mathbf{M}})\right]^{S_{ki}S_{mj}}$$



**ILLICIT**  **LEGAL**

Detectable communication pattern from Enron data. A clique surrounded by leaves keeps illegal information contained within a small group.
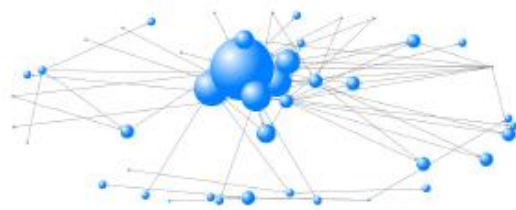


**Pattern**  **Observed Data**

Pattern matching detects exact and noisy matches to ultimately detect illegal communication patterns.
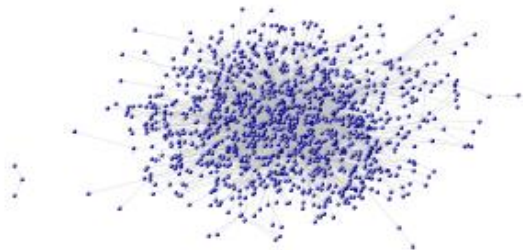
Corrupt Project 1
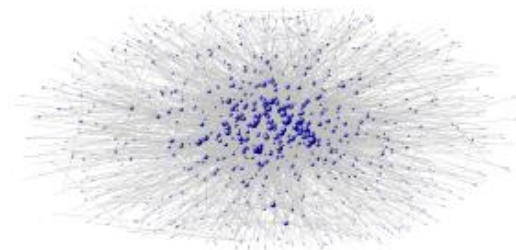
Corrupt Project 2

Corrupt Project 3

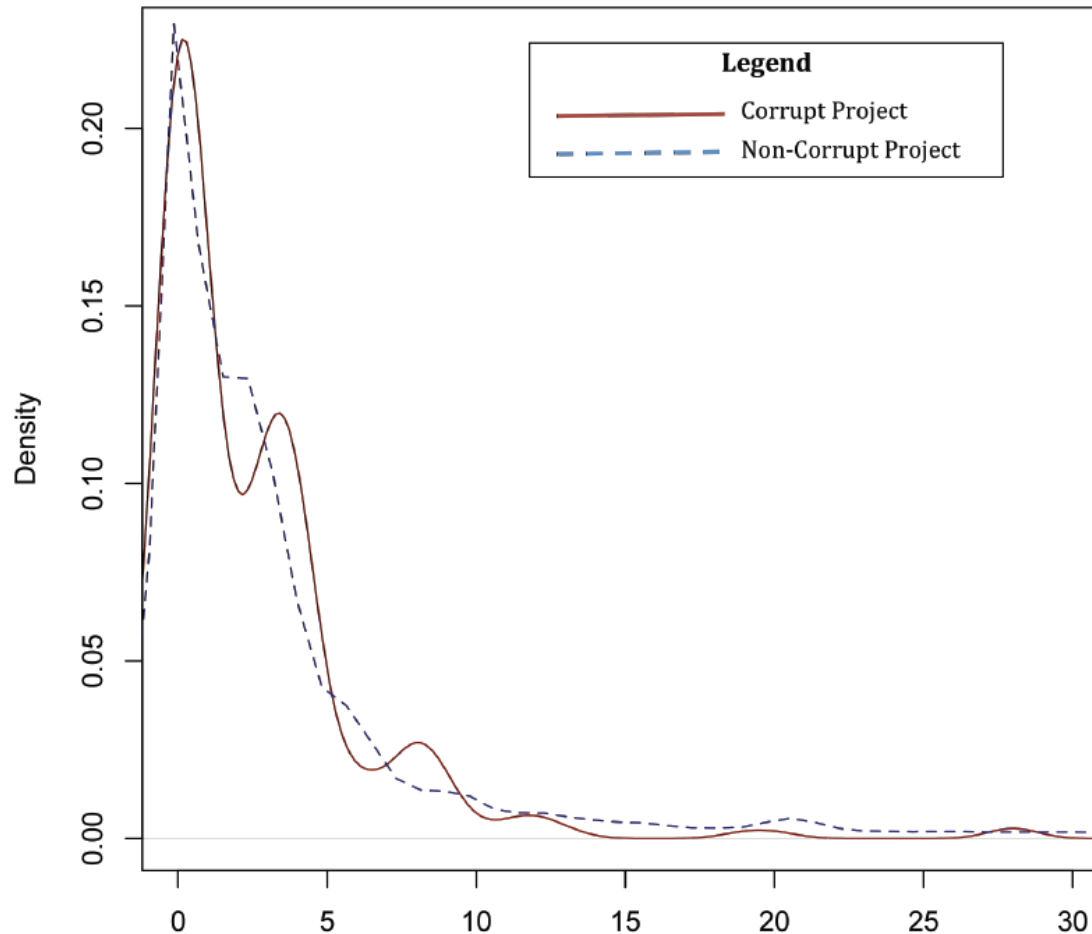Non-Corrupt Project 1

Non-Corrupt Project 2

Non-Corrupt Project 3

# Distinguishing between corrupt and noncorrupt communication requires complex network signatures

**In-degree Distribution by Project Type**



Examining a single metric such as in-degree is insufficient. Structural information appears to be required.

Given

- $S_x$ and $S_y$ independent data sources
- $P(S_y|H) \geq P(S_y)$

Then $P(H|S_x \cap S_y) \geq P(H|S_x)$.

Thus, the probability of detecting an event improves when you have more than one independent data source.
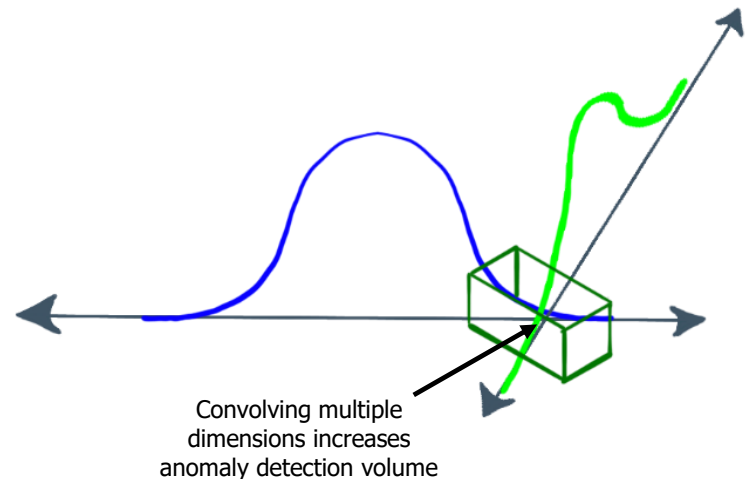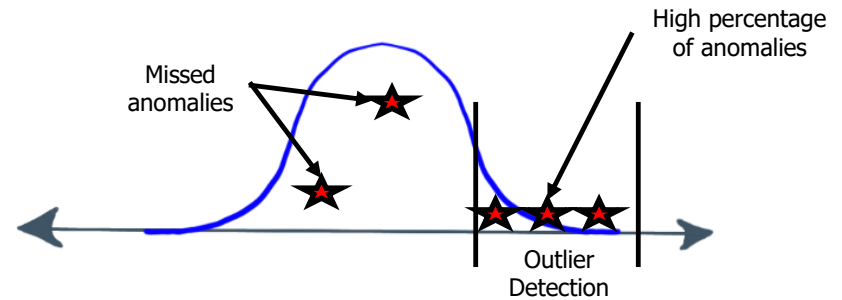
# Example: Outlier detection for anomaly identification

**Outlier detection** for identifying anomalies:

- Underlying hypothesis is that anomalies are statistical outliers along some dimension

- Commonly assume distributions are Gaussian – not necessarily true

- Low false positive rates are easily achievable AT THE EXPENSE of high false negative rates

**Insight:** Many anomalies are outliers along multiple dimensions; high false negative rates can be ameliorated:

- Develop different outlier detection algorithms for multiple dimensions

- Lower the threshold for each outlier detection algorithm → increases false positives but decreases false negatives

- Convolving the different outlier detection algorithms lowers false positives without undue impact on false negatives

- Outlier detection algorithms can be combined with more sophisticated anomaly detection techniques for further enhancement
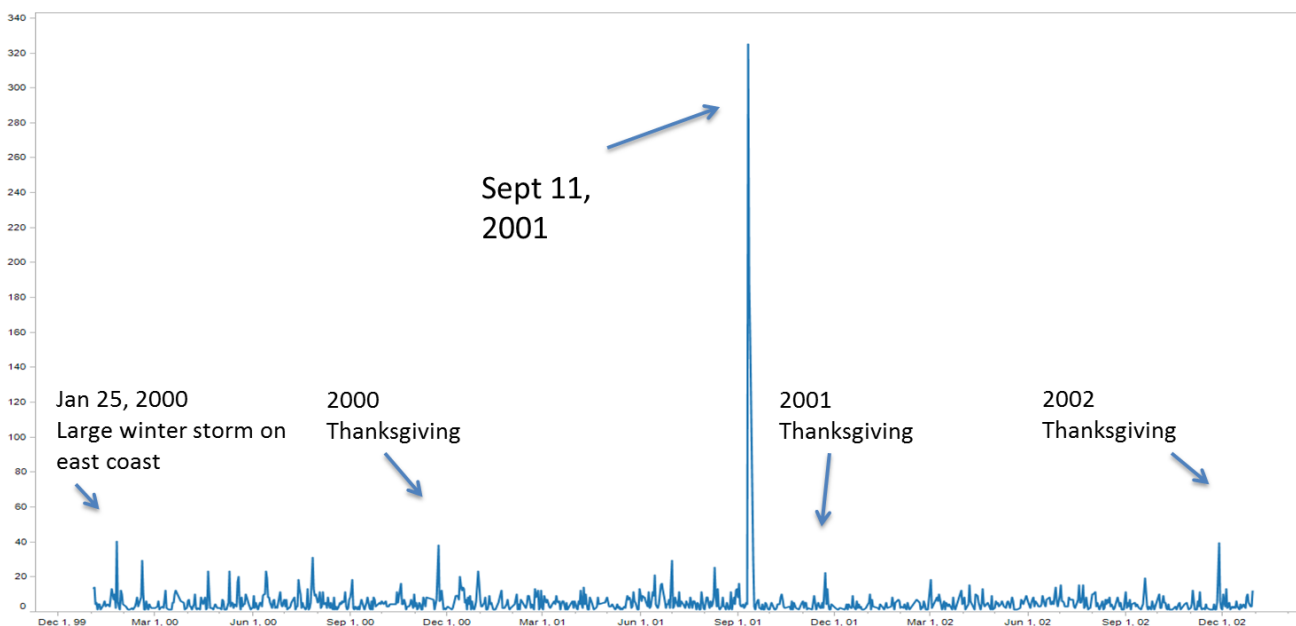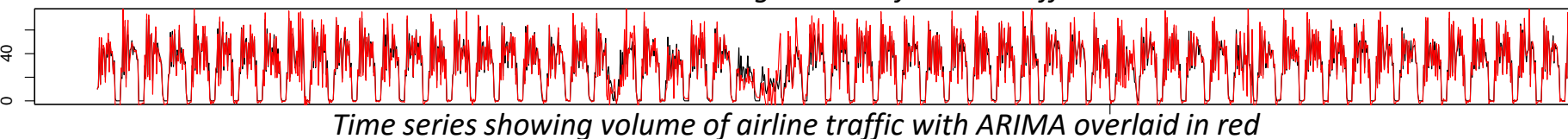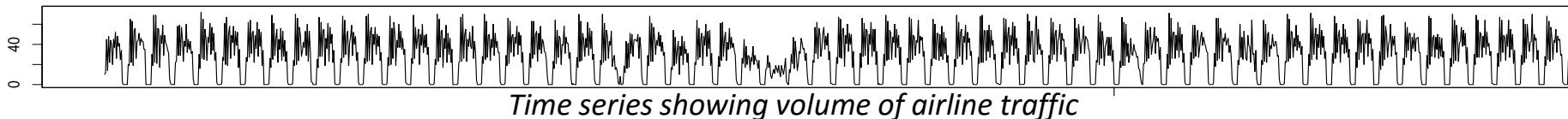
High percentage of anomalies

Missed anomalies

Outlier Detection

Convolving multiple dimensions increases anomaly detection volume

# Detect Behavior:
# Autoregressive Integrated Moving Average (ARIMA)

- Useful for behavior detection and prediction
- Incorporates seasonal trends and patterns

$$\left(1 - \sum_{i=1}^{p} \phi_i L^i\right)(1 - L)^d X_t = \left(1 + \sum_{i=1}^{q} \theta_i L^i\right)\varepsilon_t$$

*Time series showing volume of airline traffic*

*Time series showing volume of airline traffic with ARIMA overlaid in red*

Sept 11, 2001

Jan 25, 2000
Large winter storm on east coast

2000 Thanksgiving

2001 Thanksgiving

2002 Thanksgiving

*ARIMA Standardized Residuals over flight data from 2000-2002*
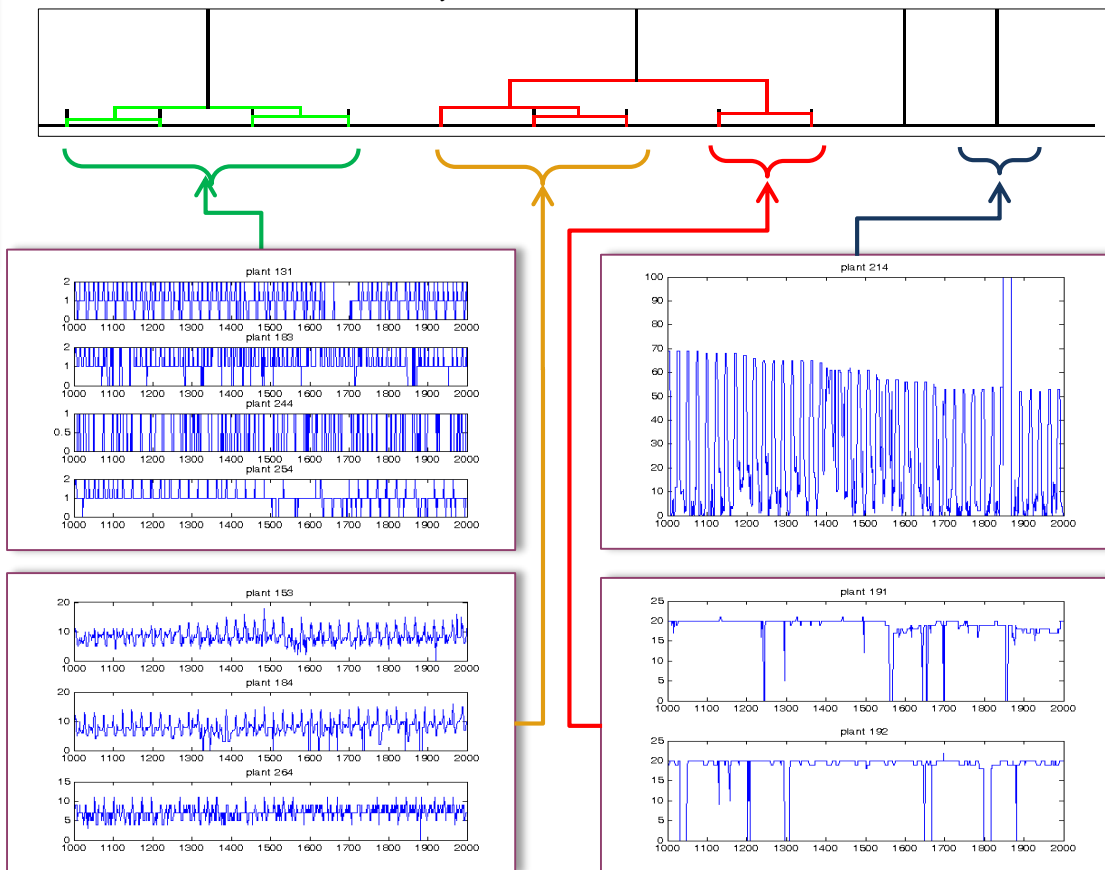
Power plant output/hour modeled as an HMM



Observation, inferred state

Predicted output and HMM state

**Clustering HMM models identifies classes of qualitatively similar signatures:**

$$d_{i,j} = D(\lambda_i, \lambda_j) = \frac{1}{T_i}\left[\log\left(P(O^i \mid \lambda_i)\right) - \log\left(P(O^i \mid \lambda_j)\right)\right]$$

$$\min_{W_1 W_2} \sum_{i=1}^{m} \left( \left\| W_2 W_1^T x^{(i)} \right\|_2^2 + \lambda \sum_{j=1}^{k} \sqrt{\epsilon + H_j \left( W_1^T x^{(i)} \right)^2} \right)$$

Edges  ⟶ Patterns in single data sources

Face Parts  ⟶ Organization parts (Teams)
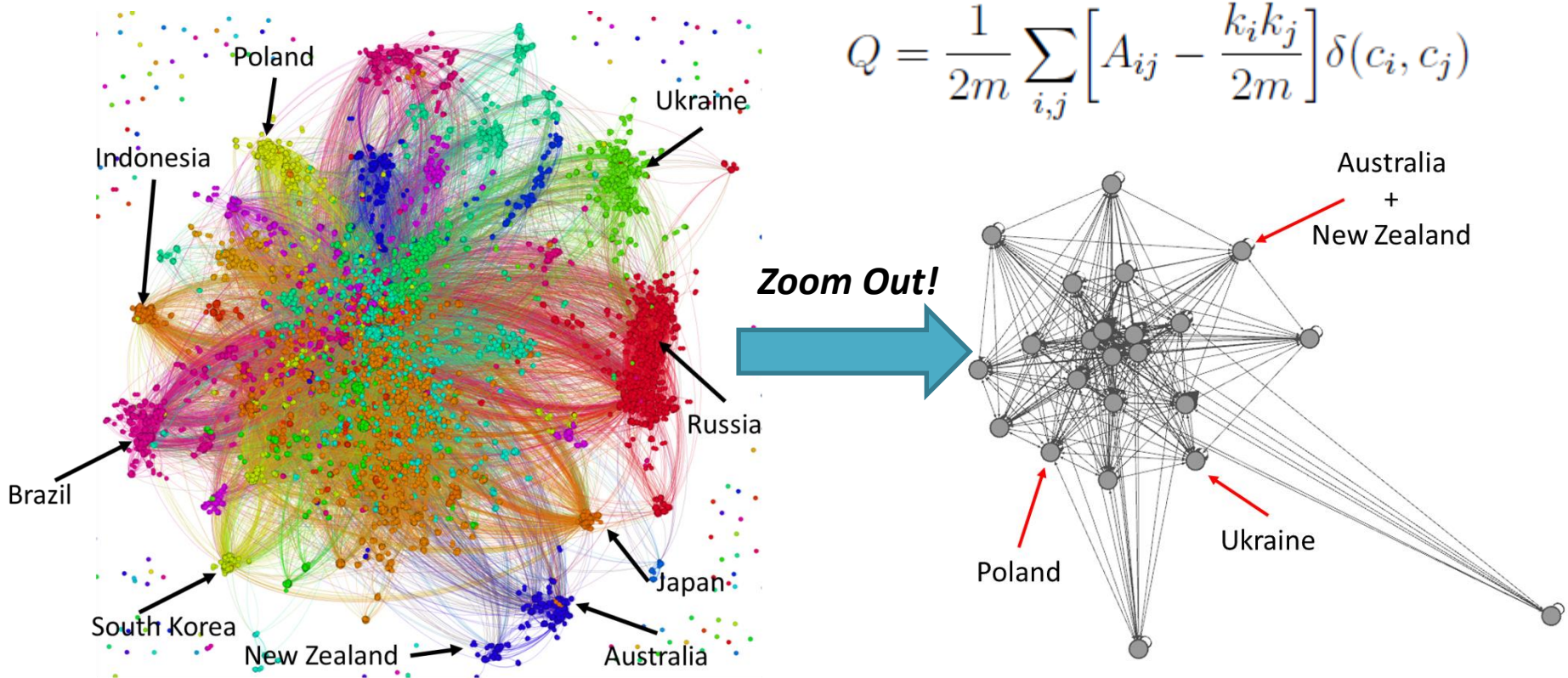
Faces  ⟶ Organizations

# Reveal Organizations: Modularity

Automated community detection based on structural characteristics

- Detects groups of nodes with high density connections amongst themselves
- Allows **hierarchical extraction of organizational** characteristics
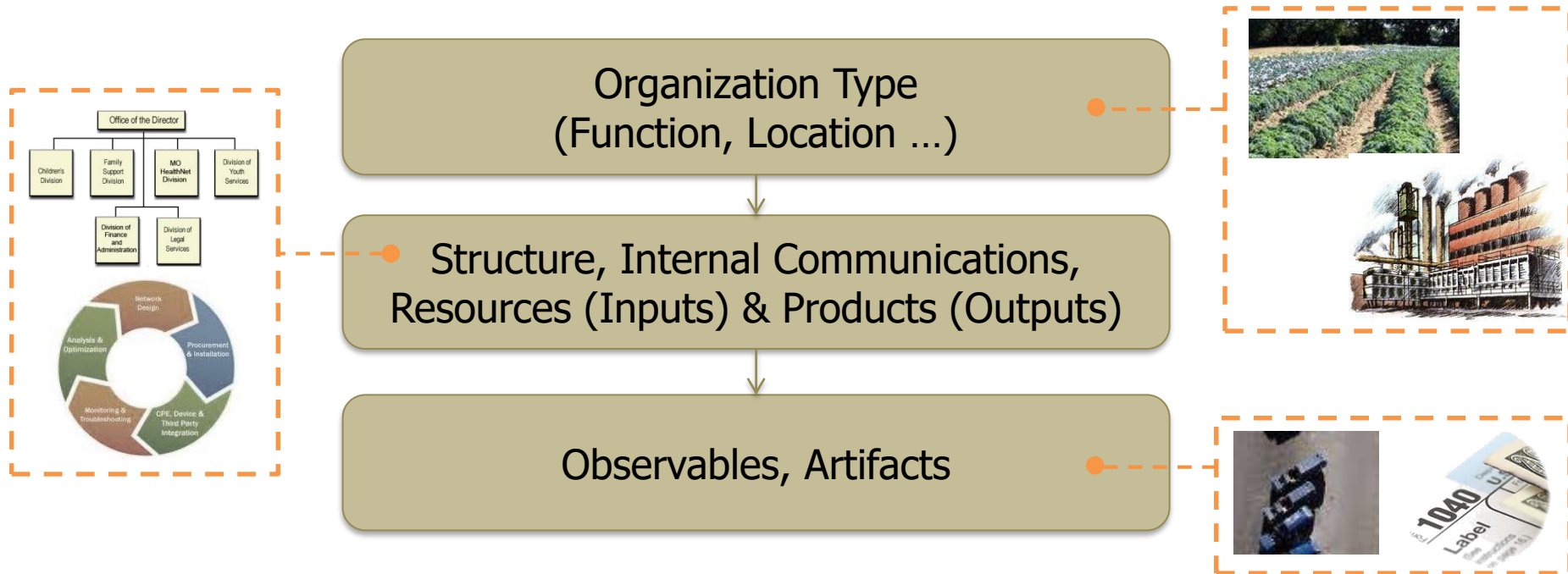
$$Q = \frac{1}{2m} \sum_{i,j} \left[ A_{ij} - \frac{k_i k_j}{2m} \right] \delta(c_i, c_j)$$



*Zoom Out!*

*Example modularity analytic shown on trace-route data. Each color shows a detected community, which ends up being closely aligned with a country / countries of interest.*

*Ability to summarize data at an organizational level and provide higher levels of abstraction / characterization*
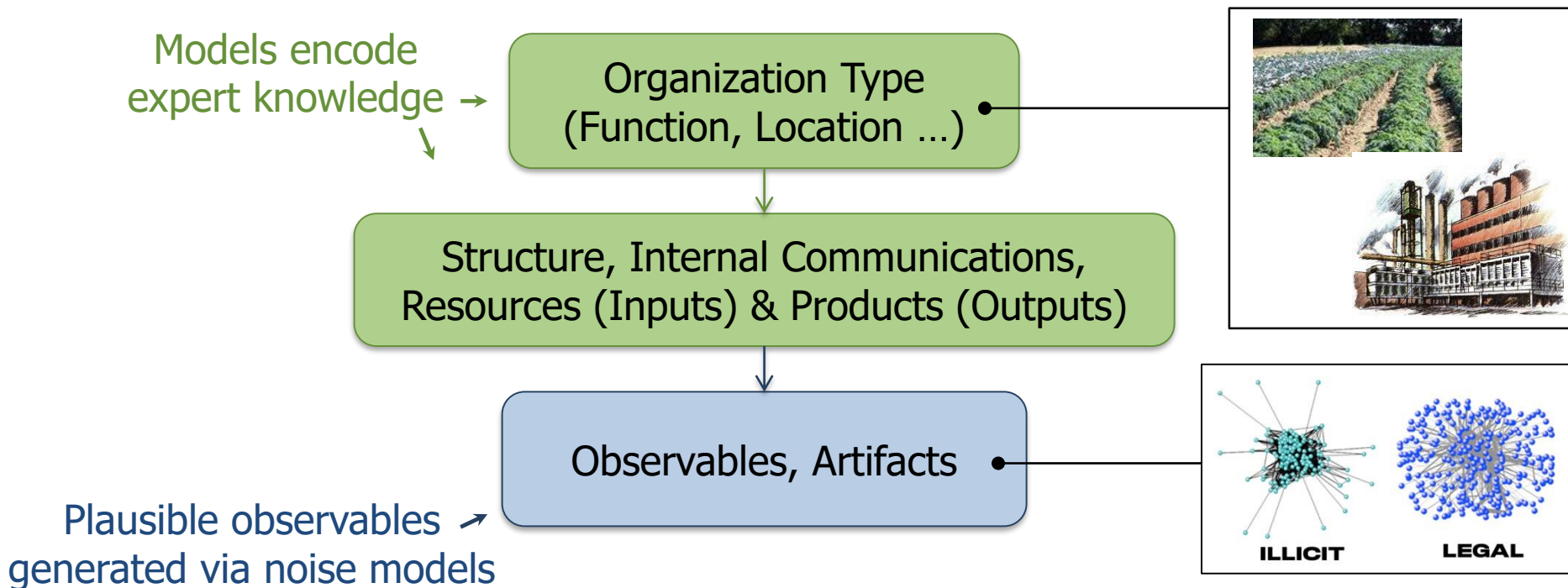
Organization Type
(Function, Location …)

Structure, Internal Communications,
Resources (Inputs) & Products (Outputs)

Observables, Artifacts

- New Technology
  - Constructing patterns of organizational activity via:
    - Generative models connecting organizations to observables
    - Network grammars constraining interactions between organizations
    - Transfer function views of organizations highlighting inputs and outputs
- Open Questions
  - What dimensions of an organization affect its structure, communication patterns, resources, and products?
  - How do the structure and processes affect the expected distributions of observables?
  - How do the expected observables and artifacts vary?

# Modeling Organizations: Deep Generative Models Connect Organizations to Observables

Models encode expert knowledge →

**Organization Type (Function, Location …)**

**Structure, Internal Communications, Resources (Inputs) & Products (Outputs)**

**Observables, Artifacts**

Plausible observables ↗ generated via noise models

ILLICIT    LEGAL

$$P(Org, Structure | Obs) \propto P(Obs | Structure) P(Structure | Org) P(Org)$$

- Organizational templates from social science provide top-down information, but observed signatures differ from theoretical predictions.
- Hierarchical generative models provide mathematics for formalizing organizational theories and connecting them to realistic signatures.

## Streaming data

Terabit/sec streaming data

## Process multi-modal datasets

ETL Processing and Refinement

Data Transformations and Filtering

$$(1 - \epsilon)C_f(S') \leq C_f(S) \leq (1 + \epsilon)C_f(S')$$

Coreset Compression

## Detect behaviors

Dynamic Tracking

Unsupervised Feature Learning

$$\left(1 - \sum_{i=1}^{p} \phi_i L^i\right)(1 - L)^d X_t = \left(1 + \sum_{i=1}^{q} \theta_i L^i\right)\varepsilon_t$$

Auto-Regressive Integrated Moving Average (ARIMA)

## Model organizations

Types → Structures → Observables

Expected & Unexpected Observables

$$P(T, S | O) \propto P(O|S)P(S|T)P(T)$$

Hierarchical generative models

## Reveal illicit organizations and relationships

Correlated behaviors

Organizational patterns

if $P(S_y | H) \geq P(S_y)$,
then $P(H | S_x \cap S_y) \geq P(H | S_x)$

Cross-modal integration improves detection