# From the Data Trenches:
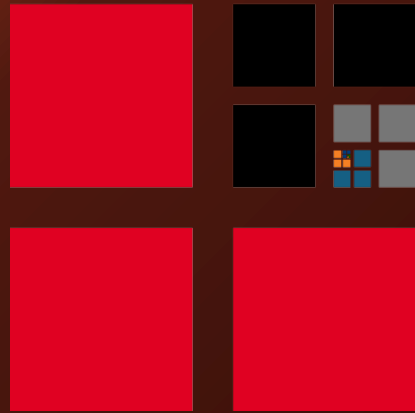## *Using Data Science for Social Good*

Eric Rozier
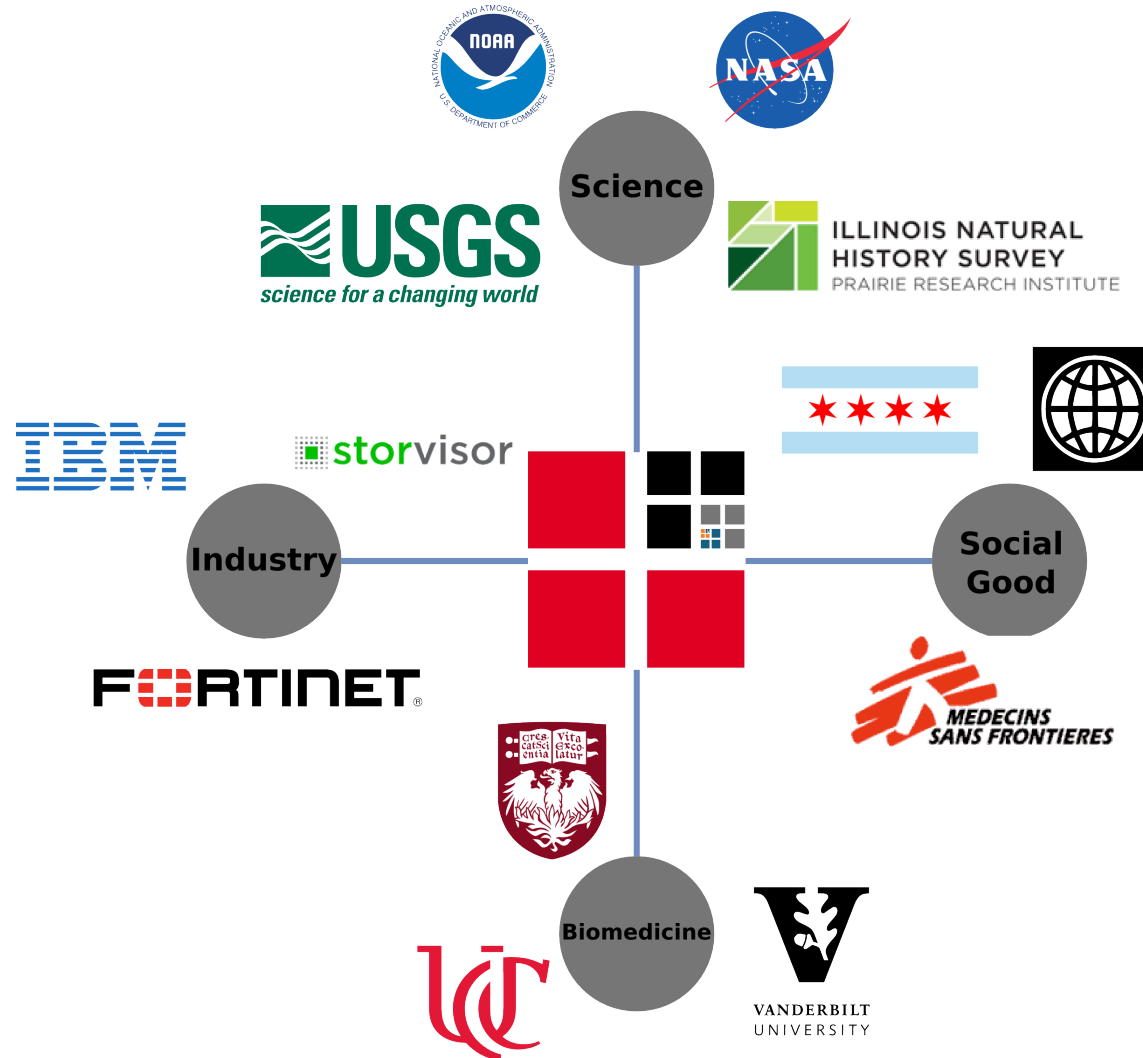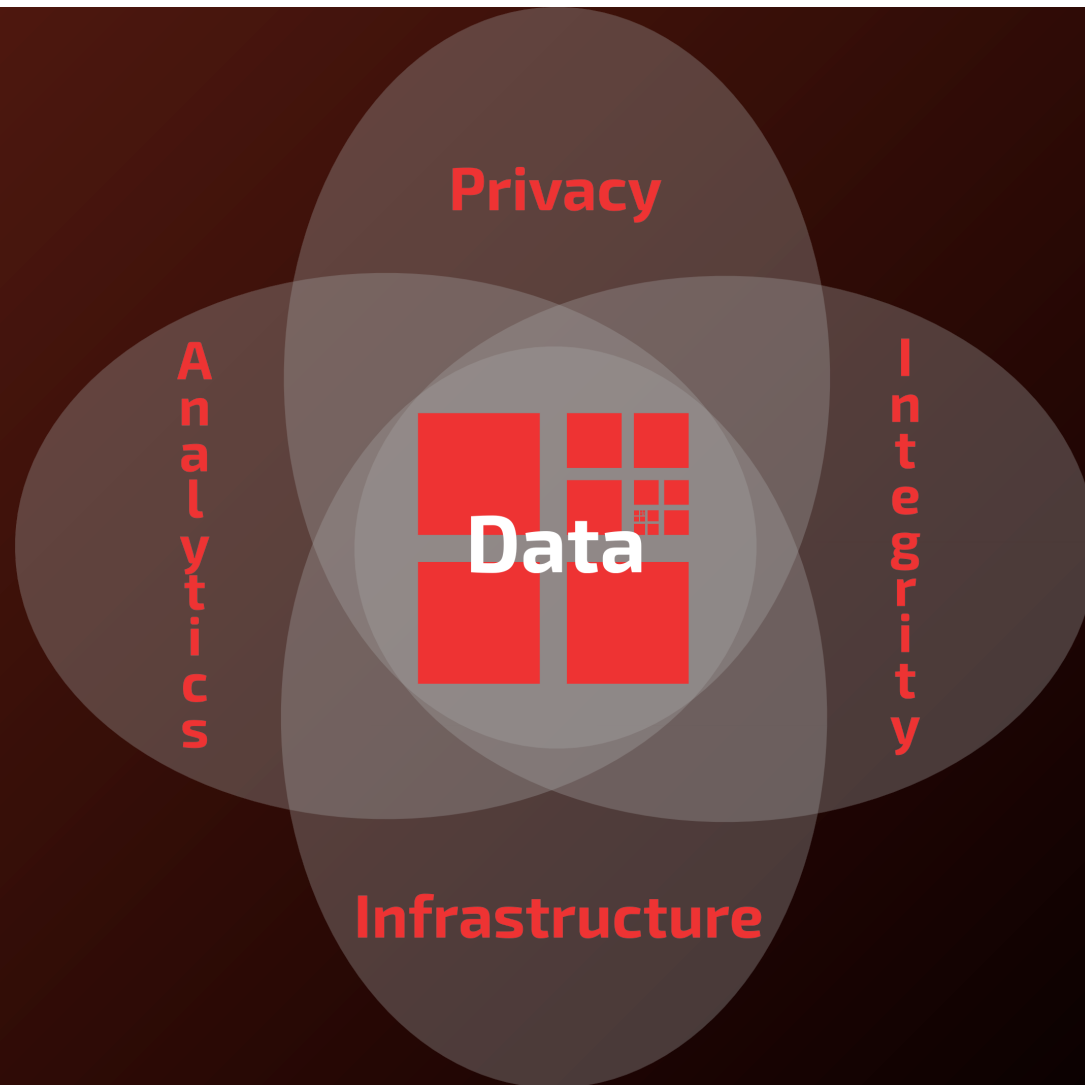
# TRUSTWORTHY
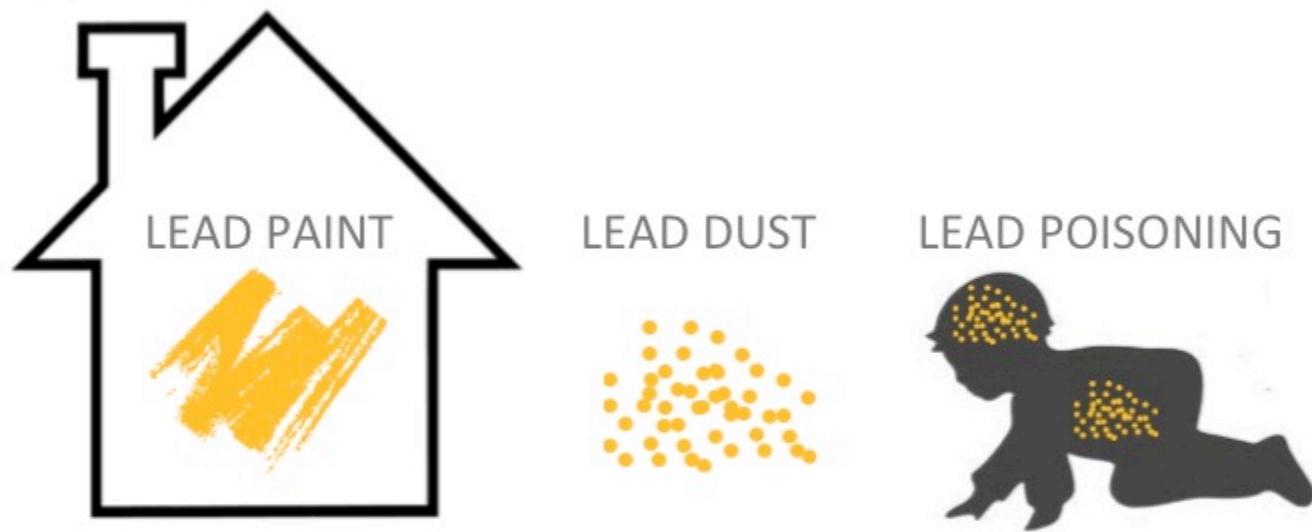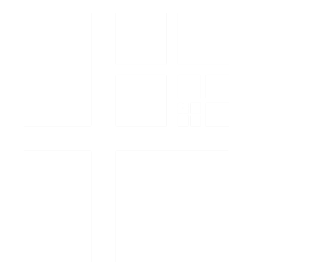## DATA ENGINEERING

# How Lead Poisoning Works

LEAD PAINT

LEAD DUST

LEAD POISONING

- NO SAFE LEVEL OF EXPOSURE
- PERMANENT HARMFUL EFFECTS

# Acting On Lead Poisoning

**Reactive**
(Current)

**Proactive**
(Target)

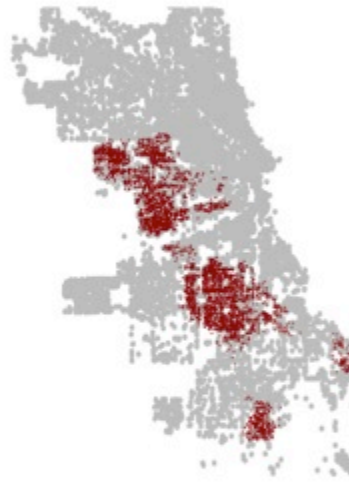| CHILD POISONED | → | HAZARD IDENTIFIED | → | HAZARD REMOVED |

# Prediction Saves Time & Money

## No Prediction



Buildings: 197,157
Time: 76 years
Money: $98 million

## Current Model



Buildings: 42,695
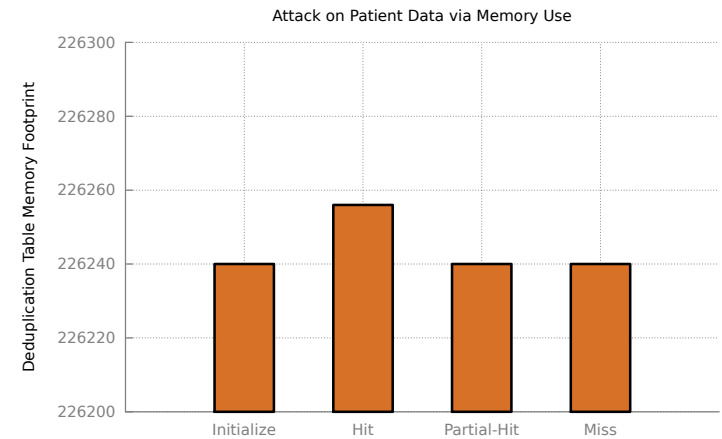Time: 16.4 years
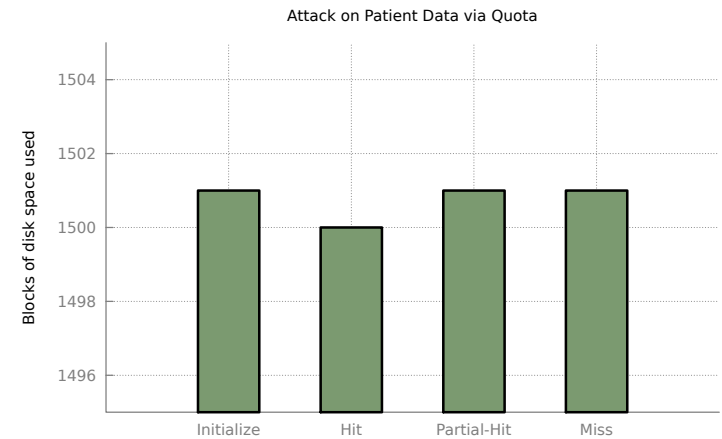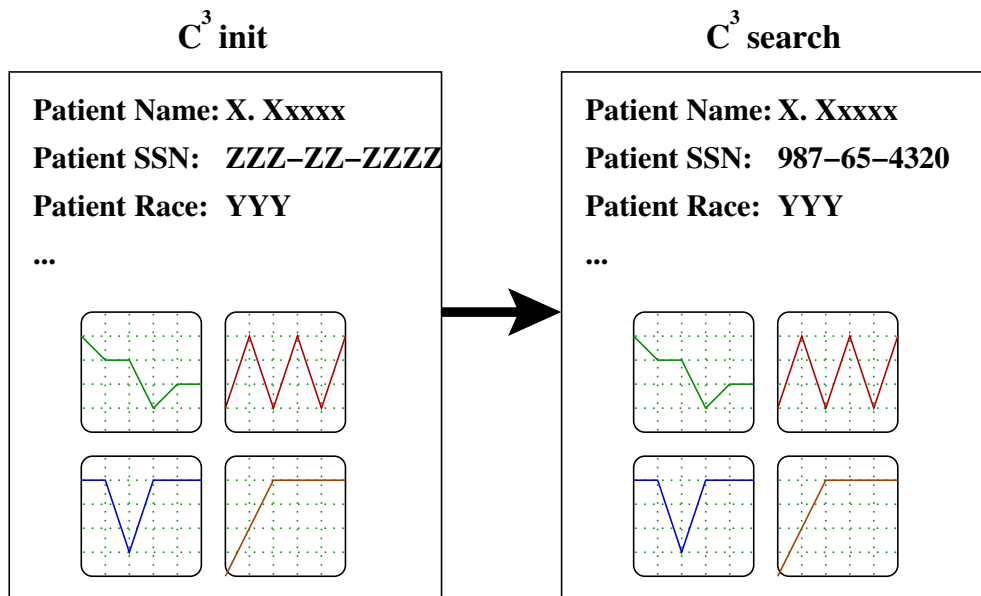Money: $21.3 million

## Model Forecast



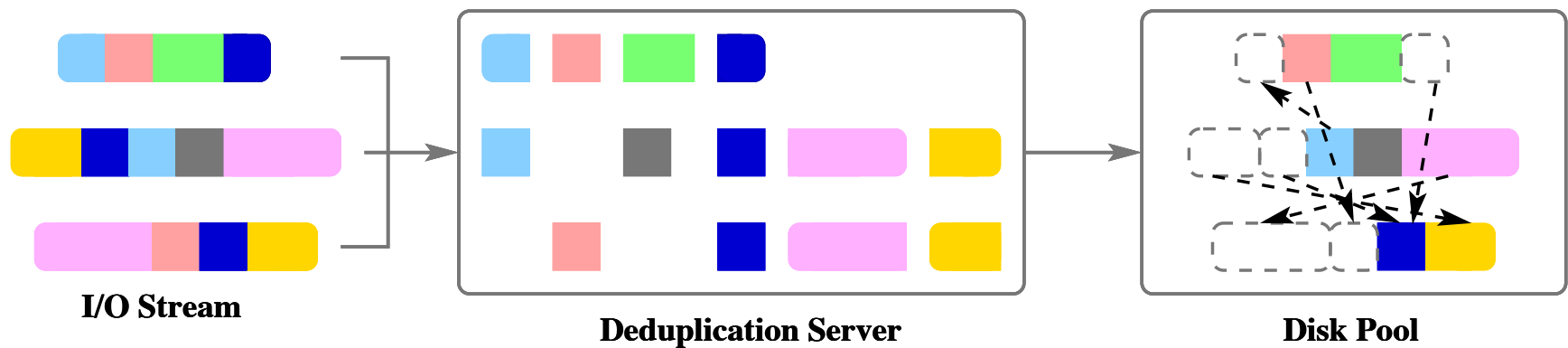Buildings: 378
Time: 2 months
Money: $189,000

# Two Problems

- Systems and data which are not protected

- Systems and data with unknown vulnerabilities

# Server-Side Dedup

### C³ init

**Patient Name: X. Xxxxx**

**Patient SSN:   ZZZ–ZZ–ZZZZ**

**Patient Race:  YYY**

**...**

### C³ search

**Patient Name: X. Xxxxx**

**Patient SSN:   987–65–4320**

**Patient Race:  YYY**

**...**

Attack on Patient Data via Quota

Attack on Patient Data via Memory Use

# Understanding Deduplication



I/O Stream

Deduplication Server

Disk Pool

# Benefit of Cross Deduplication?



Deduplication ratio EMR patient data

Deduplication ratio of sea ice records for the Southern Hemisphere during 2007 stored in the cloud on a daily basis

# Dealing with Untrusted Storage

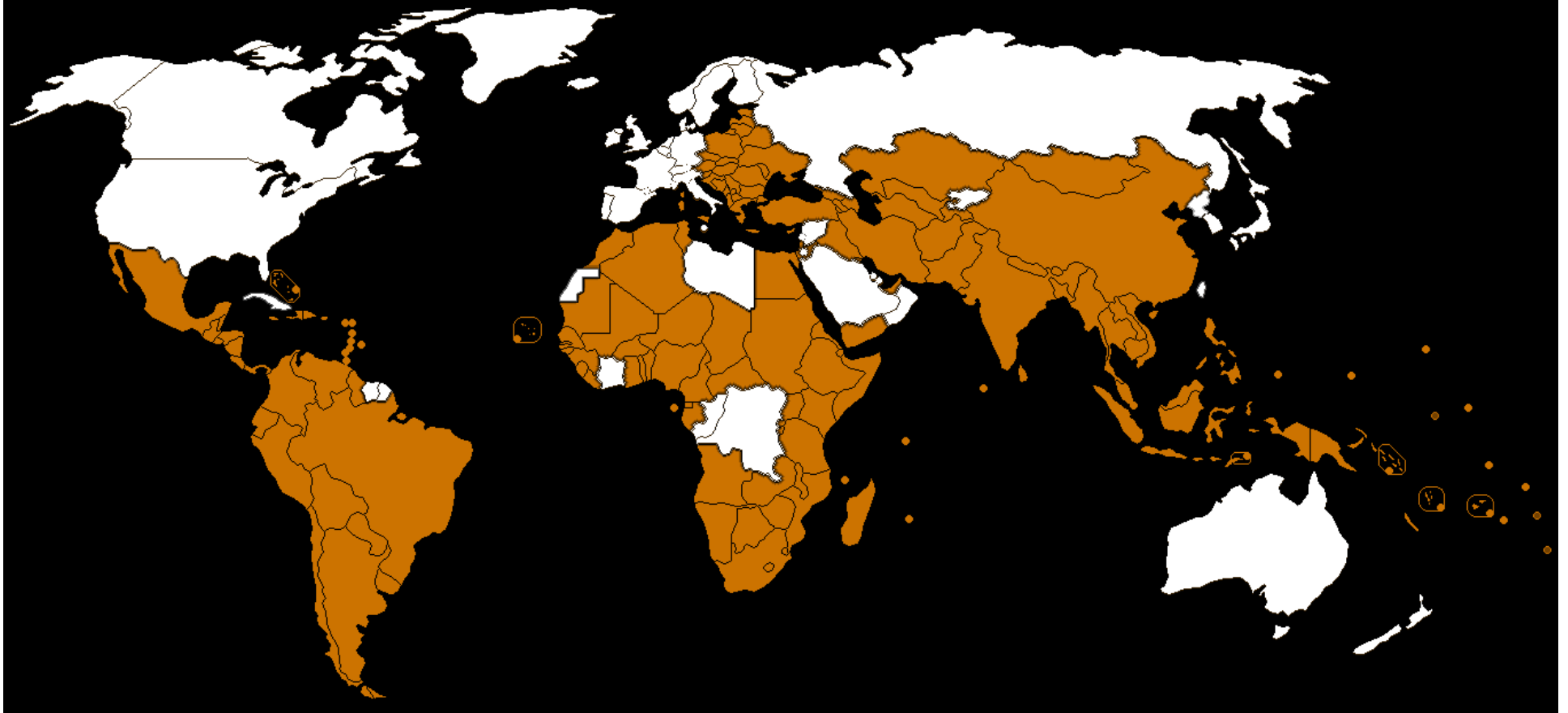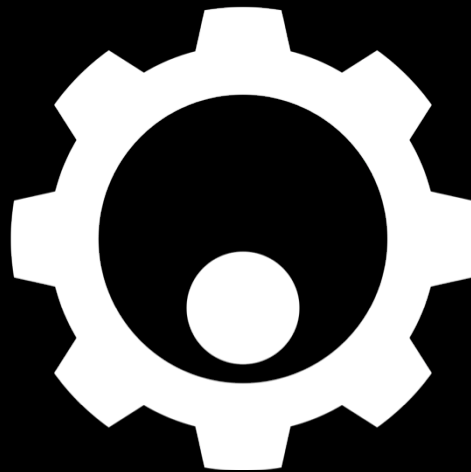# International Procurement
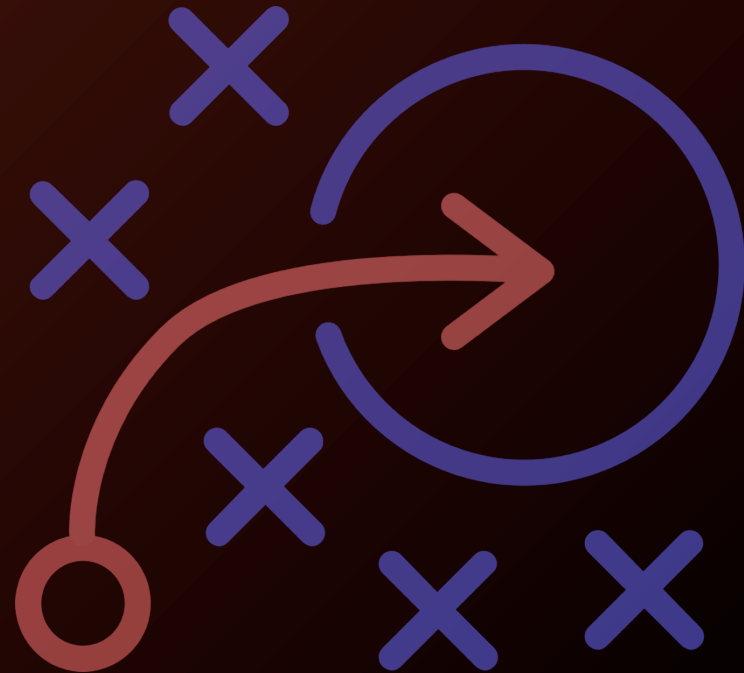
International Procurement

Automated Reasoning and Machine Intelligence

# Strategy of Data Integrity Attacks

- Data pollution
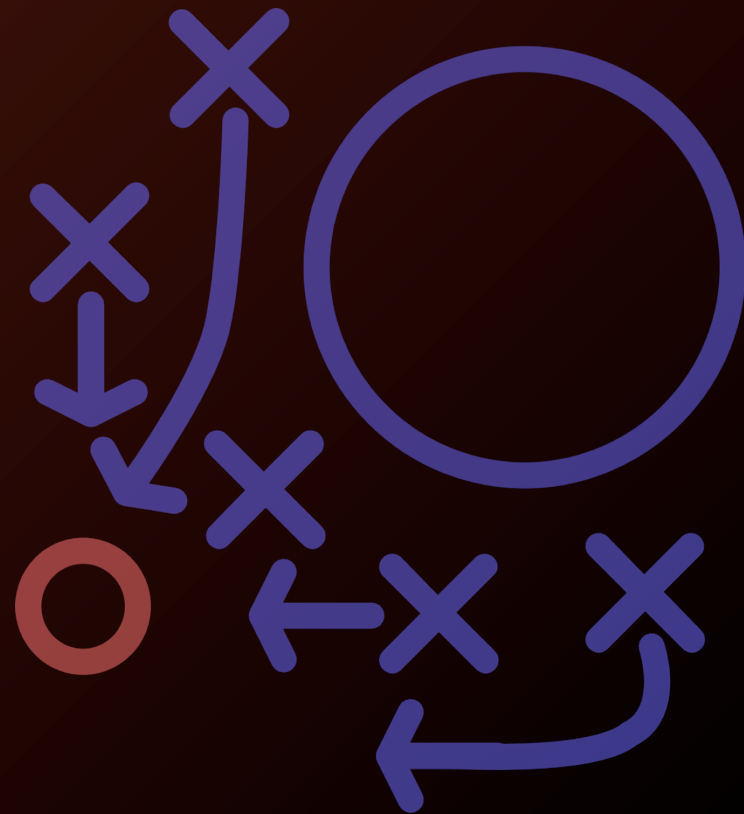- Data falsification
- Data blending

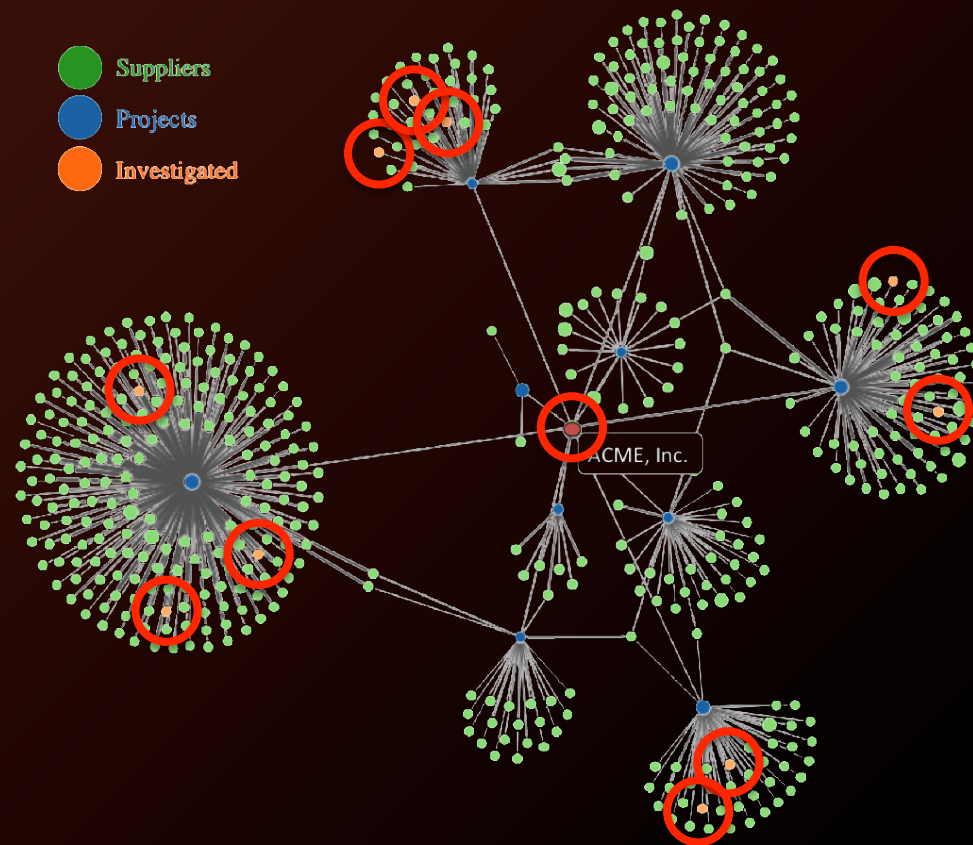Tools Designed for Non-Malicious Environments
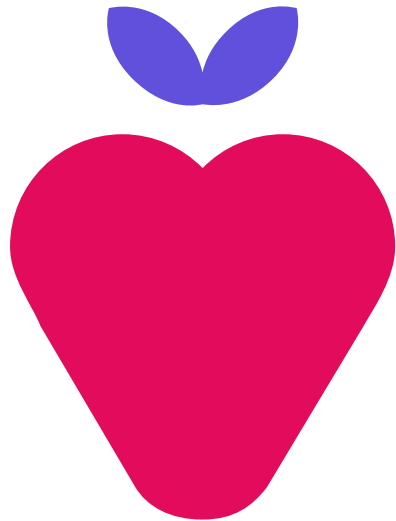
# Strategy of Data Integrity Defense

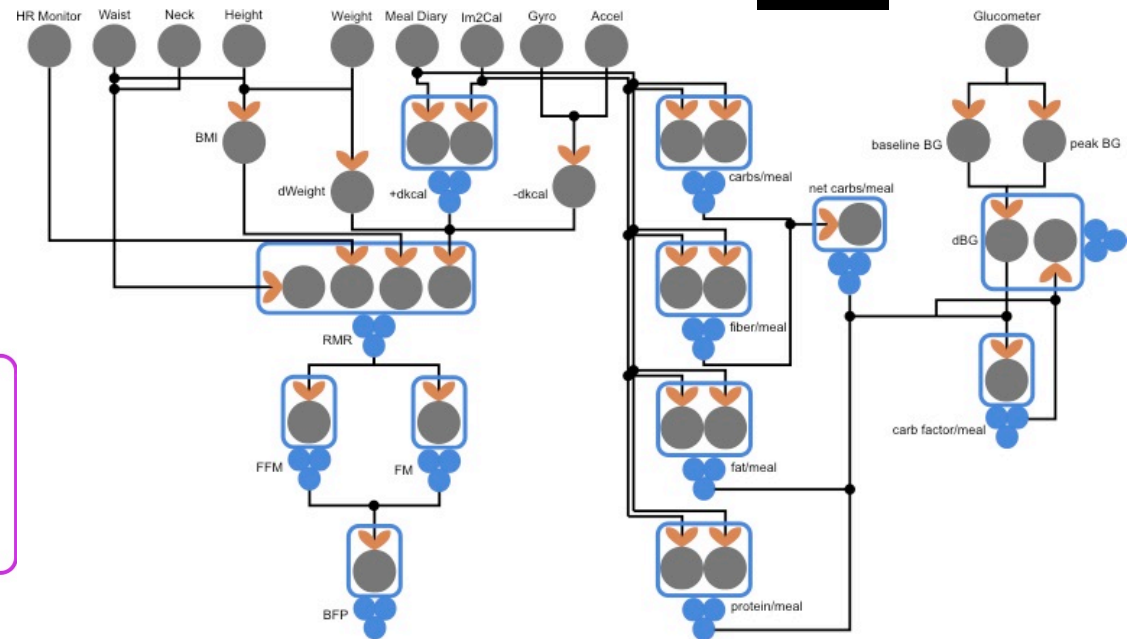- Syntactic Algorithms
- Semantic Algorithms
- Data Superiority

# Fooled by Statistics

High Risk

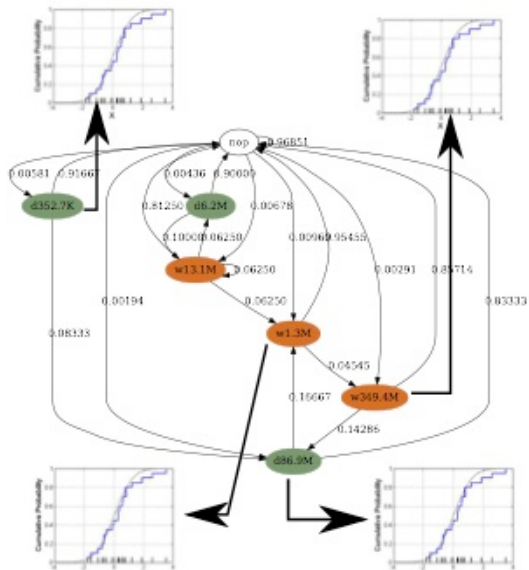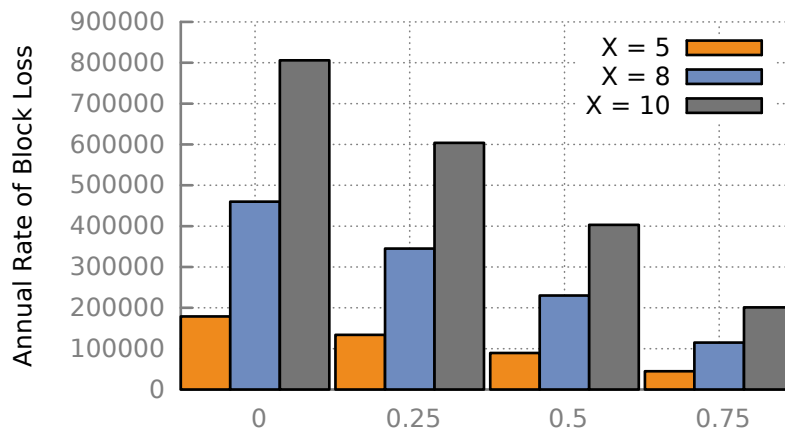Low Risk

Suppliers

Projects

Investigated

ACME, Inc.

Project Fraise

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| HR Monitor | Waist | Neck | Height | Weight | Meal Diary | Im2Cal | Gyro | Accel | | Glucometer |

BMI

dWeight   +dkcal   -dkcal

carbs/meal

net carbs/meal

baseline BG   peak BG

RMR

fiber/meal

dBG

FFM   FM

fat/meal

carb factor/meal

BFP

protein/meal

Privacy
Preserving
Searchable
Encryption

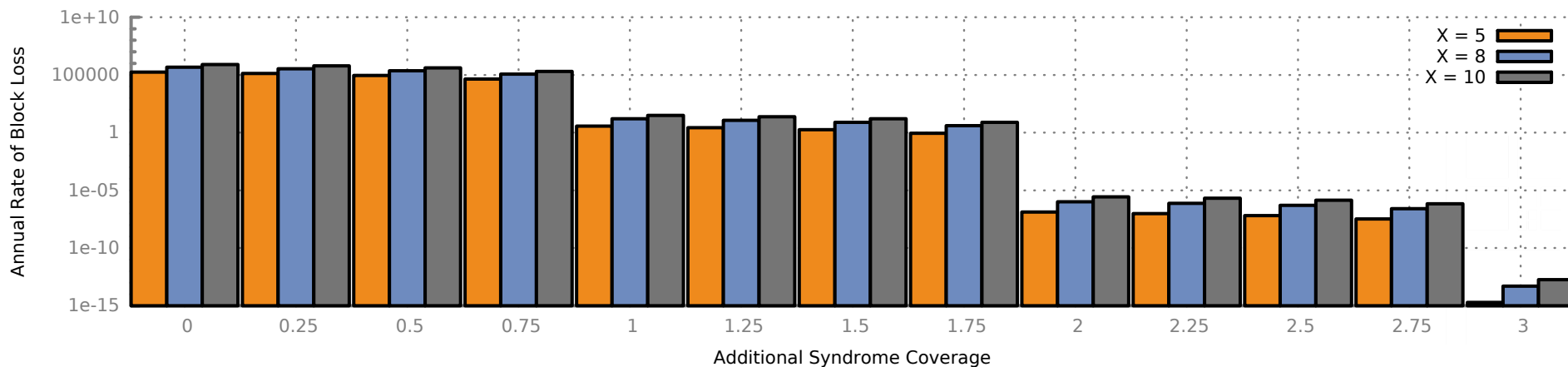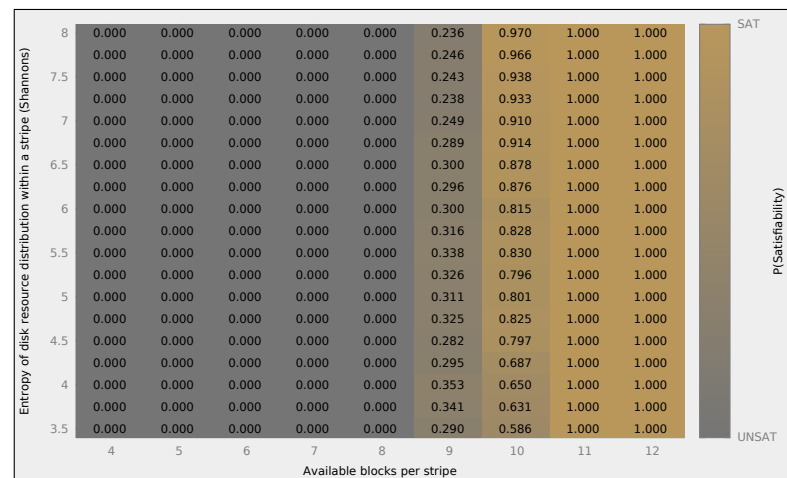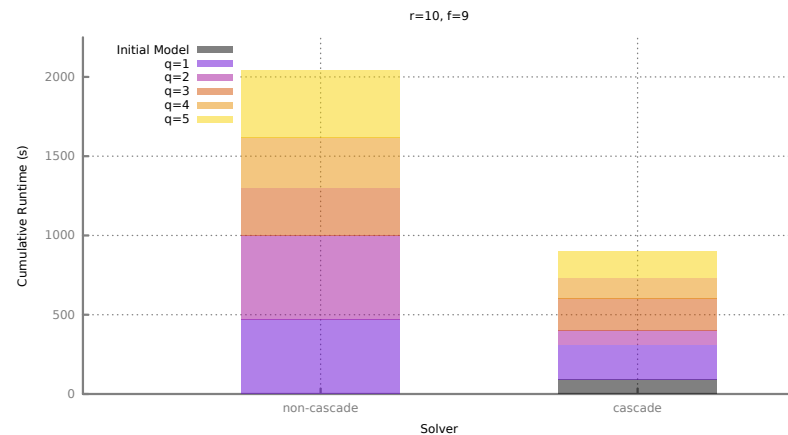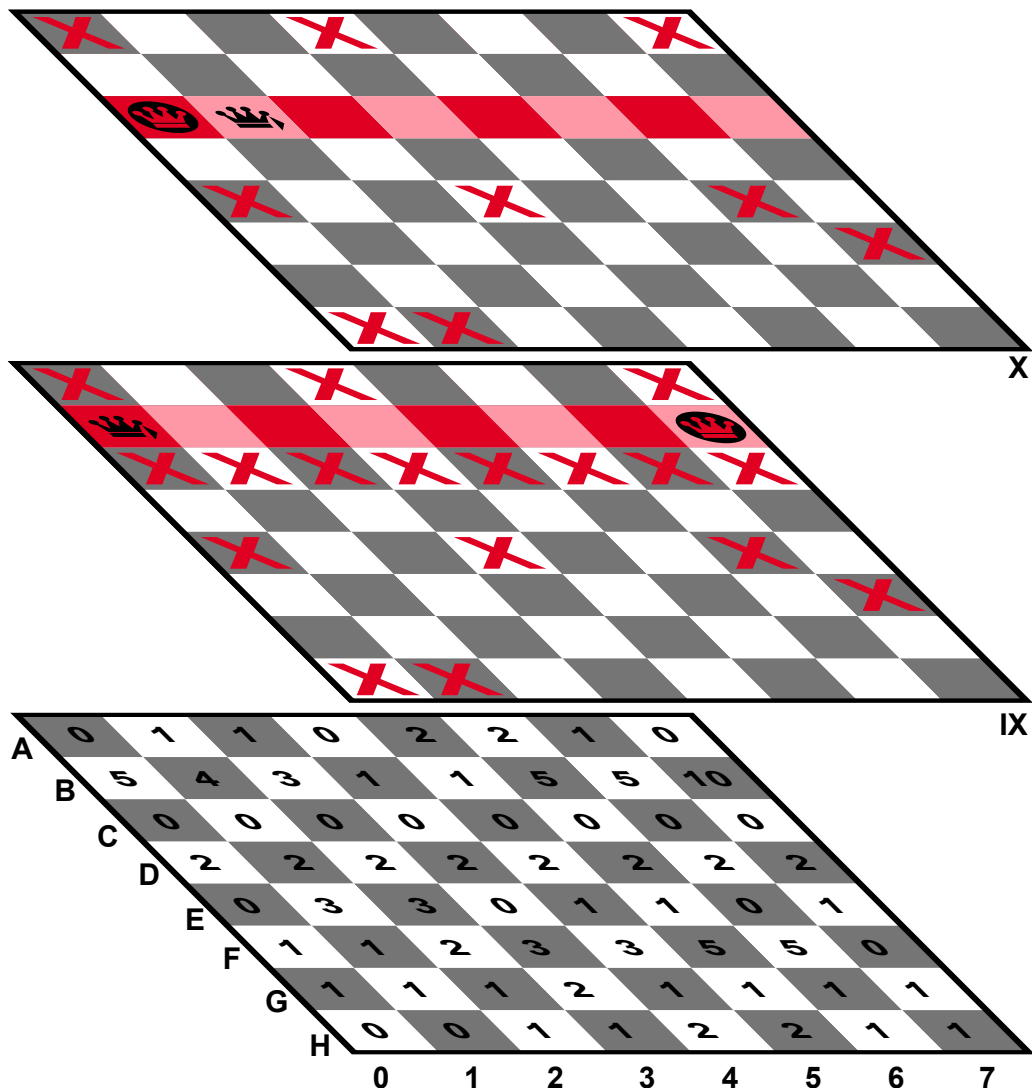| Data Aquisition Layer |
| Vineyard Data Engineering Layer |
| Health Indicators |

Annual Rate of Block Loss for a 1 Petabyte storage system with X+1 initial RAID configuration.



Annual Rate of Block Loss for a 1 Petabyte storage system with X+1 initial RAID configuration.

- How do we protect data?

- How do we ensure data integrity?

- How do we engineer systems for the new adversarial environment?