

# Distributed Summarization of Dynamic Data

Emma Alexander and Eric Balkanski

CS 262 Project Presentation  
April 25<sup>th</sup>, 2016




# SUMMARIZING CRIME DATA

**Dataset**

<u>City</u>	<u>Type</u>	<u>Time</u>	<u>Arrest</u>
NY	Assault	Night	Yes
NY	Theft	Night	Yes
NY	Assault	Night	No
LA	Theft	Morning	Yes
LA	Theft	Evening	No

**Summary**



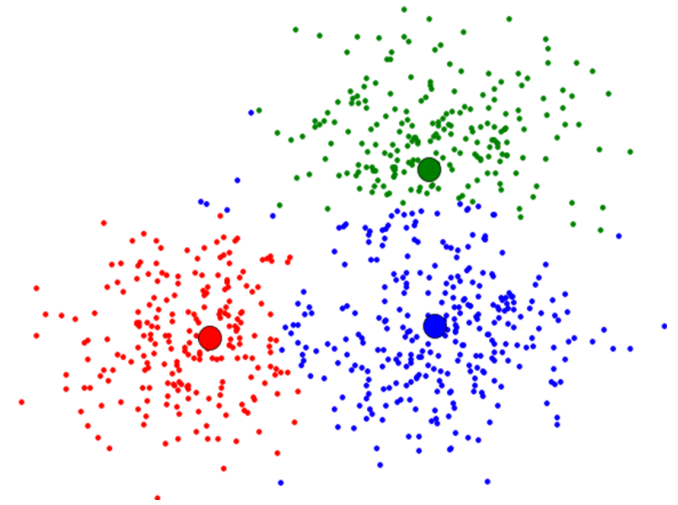
NY	Assault	Night	Yes
LA	Theft	Morning	Yes



# DATA SUMMARIZATION AS OPTIMIZATION

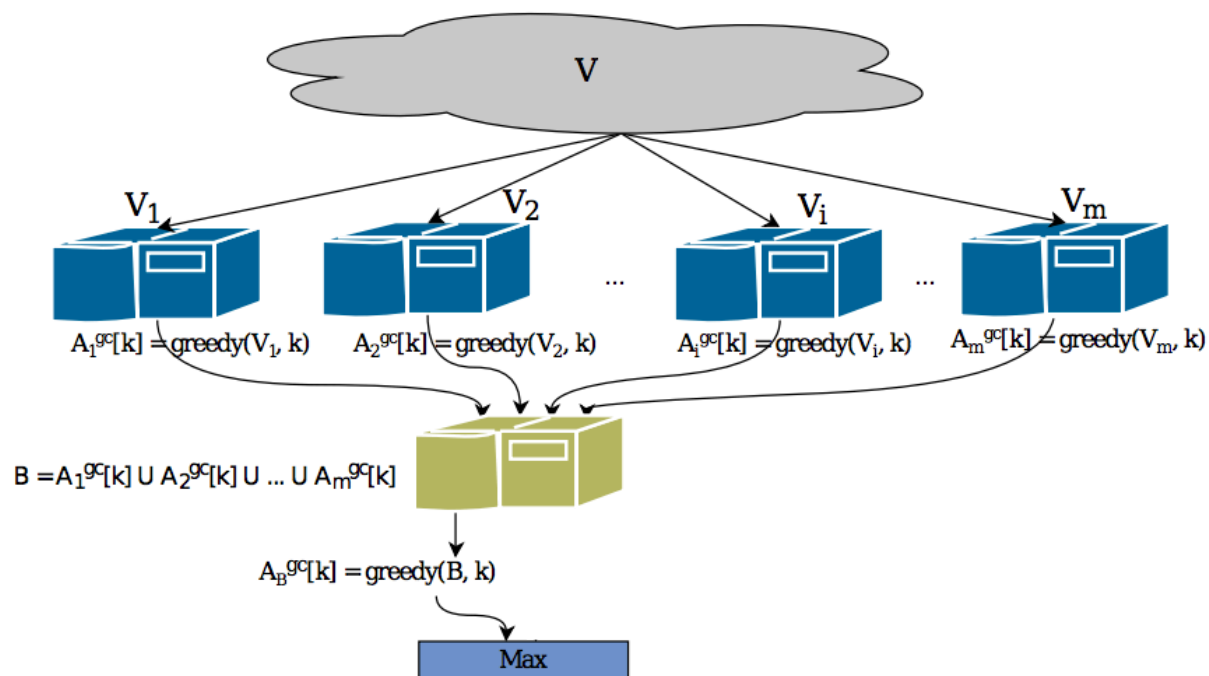
- Input:
  - **set of elements**  $V$  (e.g., set of crimes)
  - **function**  $f(S)$  measures how well  $S$  summarizes  $V$ 
    - via **clustering**
- Optimization problem:

$$\max_{\substack{|S| \leq k \\ S \subseteq V}} f(S)$$







# DISTRIBUTED DATA SUMMARIZATION

- What if  $V$  is too large to: to efficiently compute summary? fit on single machine?
- Distributed approach: [Mirzasoleiman et al. '13 and '15]



# TRADEOFFS

As number of machines  $m \nearrow \dots$

- Runtime 
- Memory per machine 
- Approximation ratio  $\alpha = f(A^*) / f(A)$   
(quality of  $A$  compared to optimal solution  $A^*$ ) 
- Communication complexity 



# OUR CONTRIBUTION

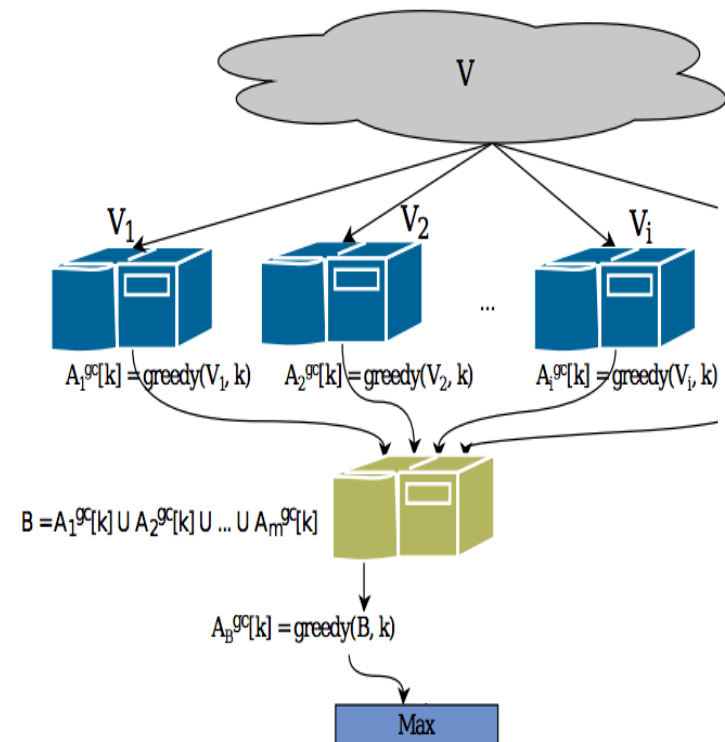
- Real world datasets are **dynamic** with
  - **Insertions** (e.g., new crimes happen everyday)
  - **Deletions** (e.g., erroneous or outdated records)
- Naïve solution: rerun entire algorithm for every insertion and deletion
- We can achieve **better communication complexity** ...





# OUR SOLUTION

- Insertion or deletion for  $V_i$
- $A_i^{old} :=$  last set sent from machine  $i$

1. Update local solution  $A_i$
2. If  $d(A_i, A_i^{old}) \geq t$ :
  - Send  $A_i$  to central machine
3. Update central solution  $A$  when receive updated  $A_i$



# TRADEOFFS

- Trigger for sending a new local solution:  $d(A_i, A_i^{old}) \geq t$
- As threshold  $t$  ↗ :
  - Approximation ratio 
  - Communication complexity 





# FAILURE

- On local machine, either
  1. central machine waits for  $m - T$  local solutions  $A_i$ 
    - $T$  fault tolerance
    - Approximation ratio
  - $T+1$  replicas of each element on different machines
    - $T$  fault tolerance
    - Memory per machine
- On central machine, either
  - leader election
  - $T+1$  replicas of central machine

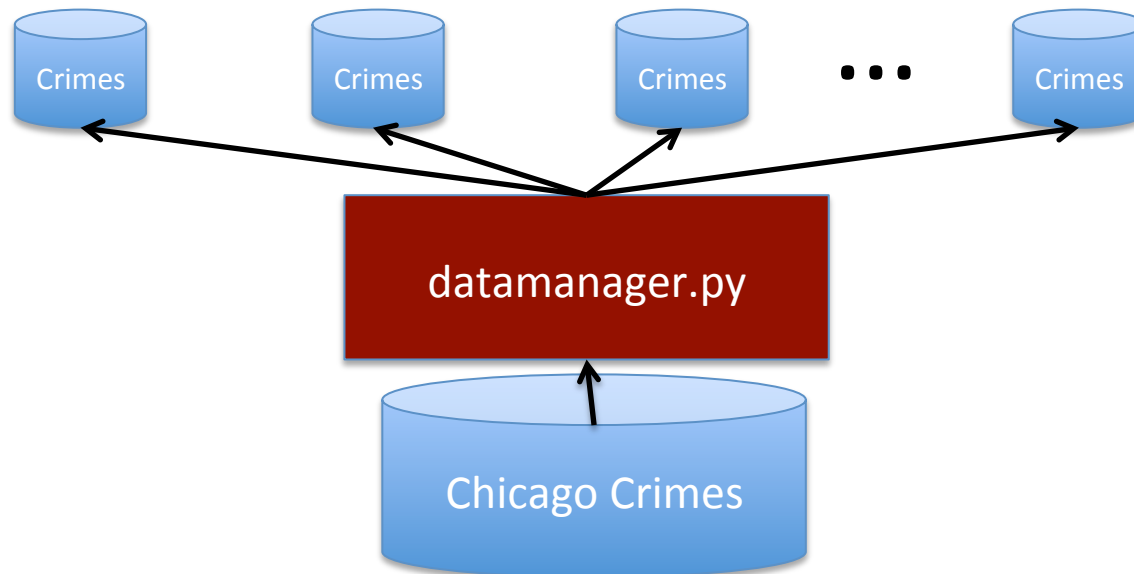


# IMPLEMENTATION...

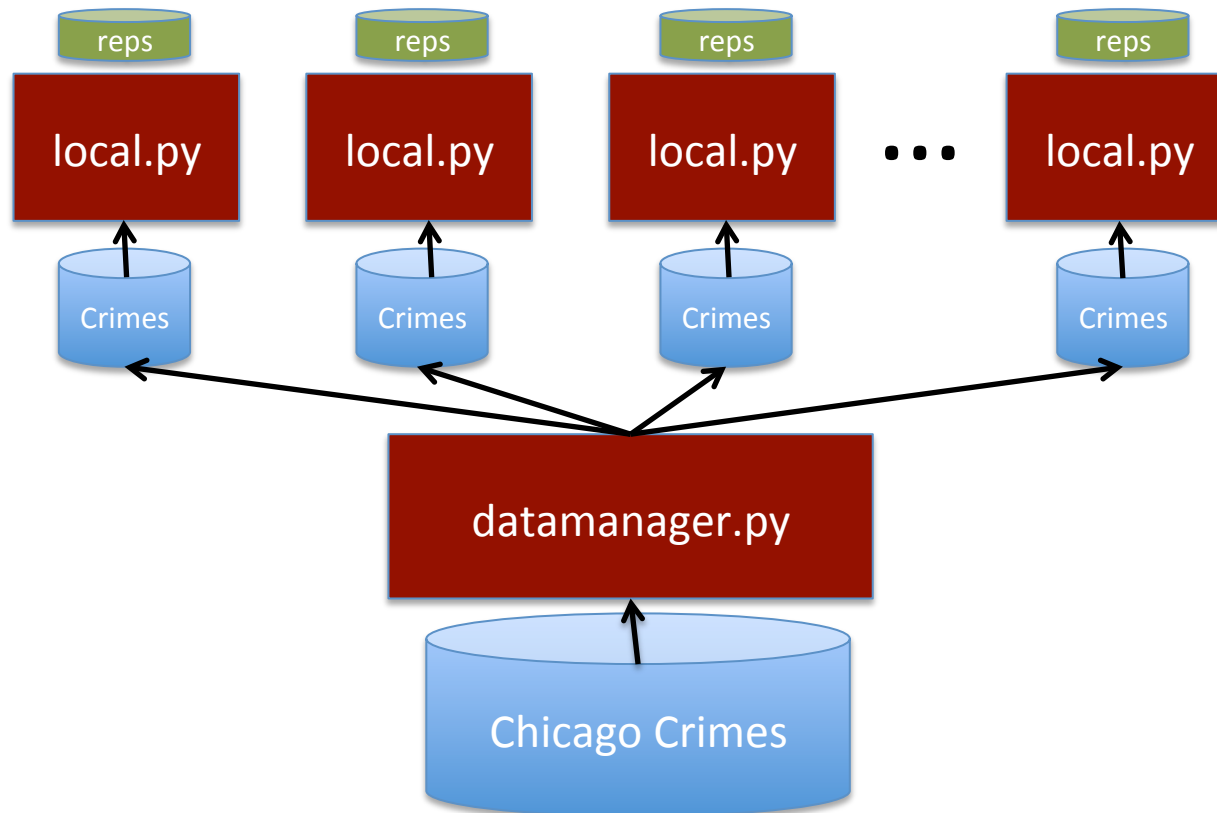




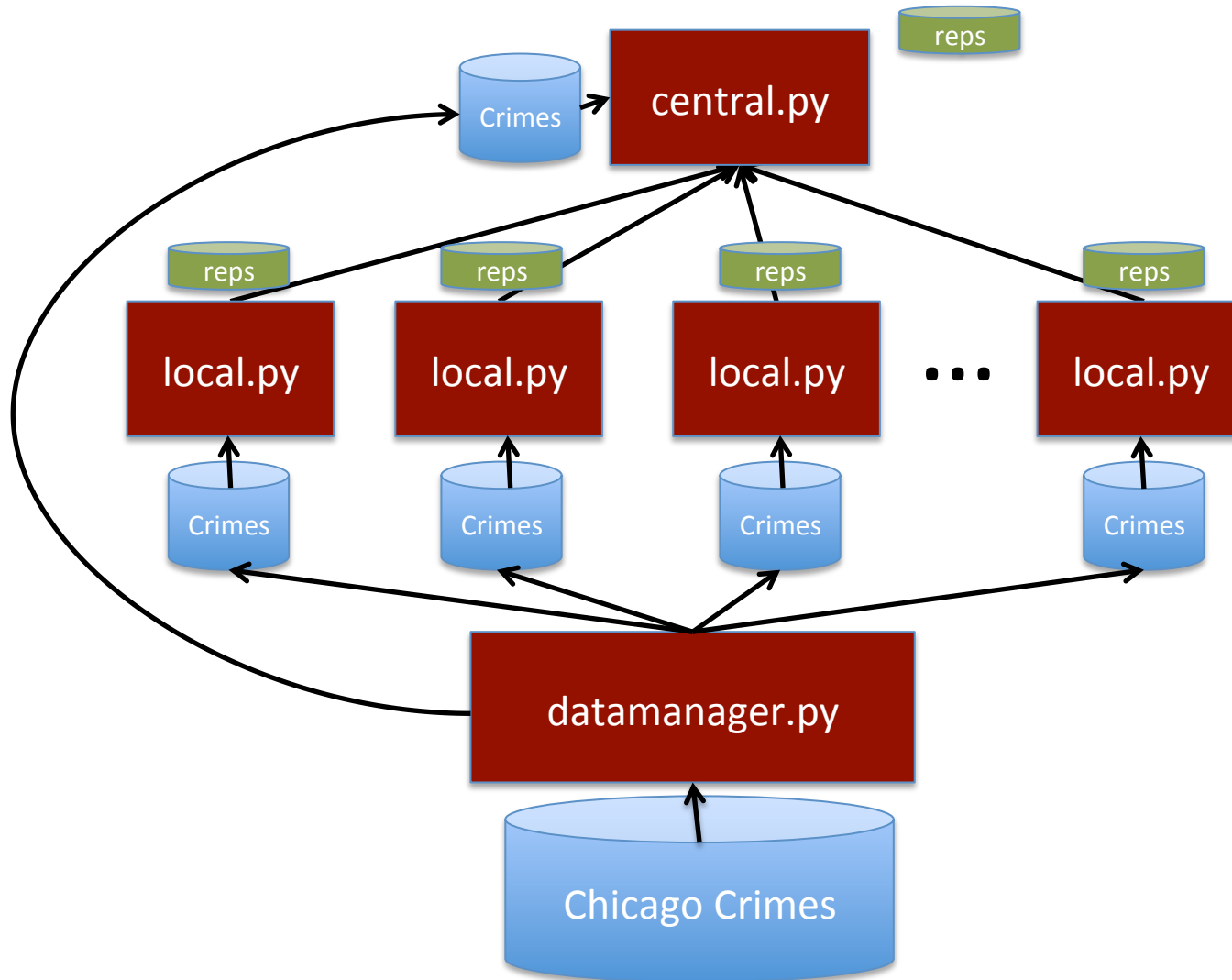
# SIMULATION



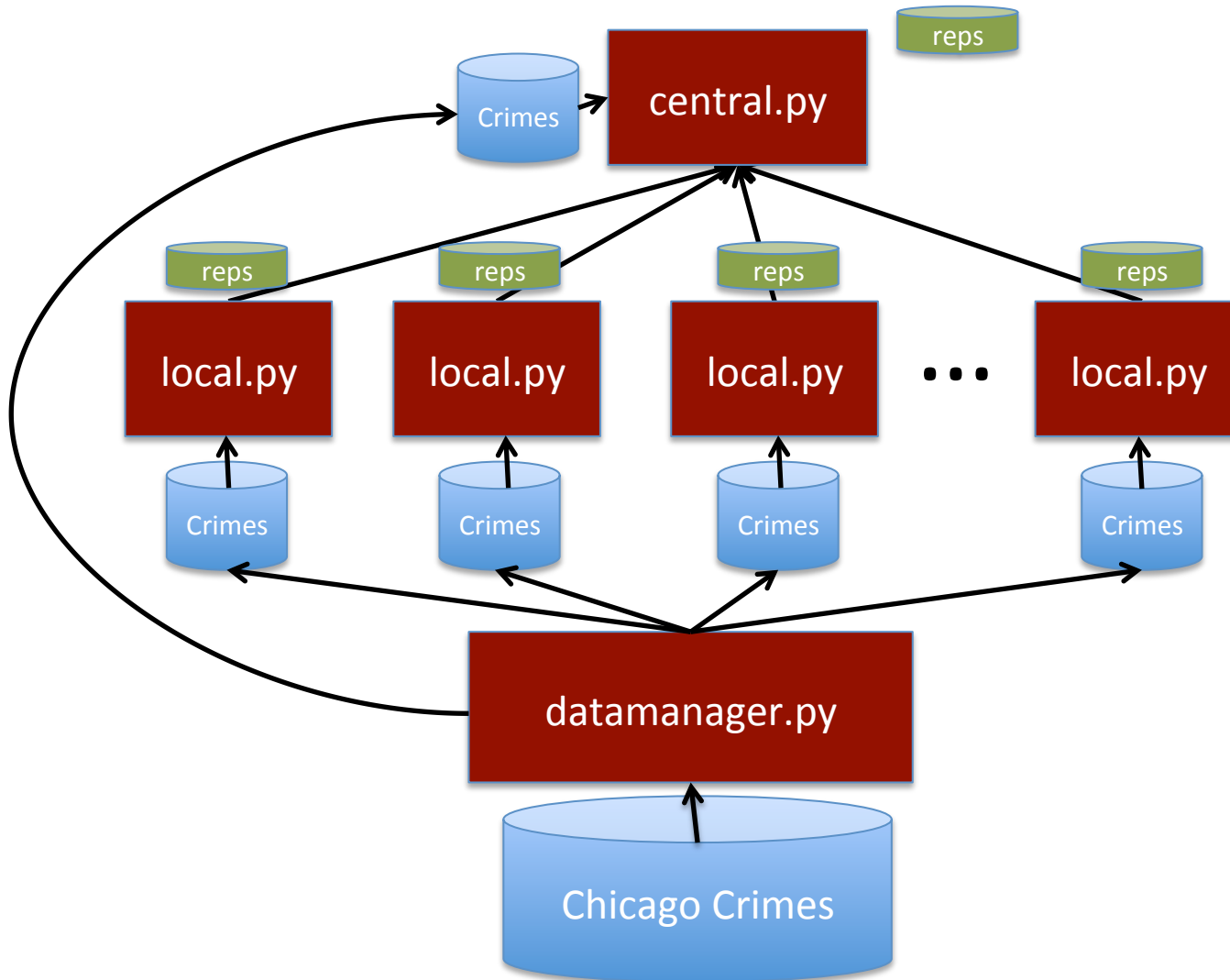
# SIMULATION



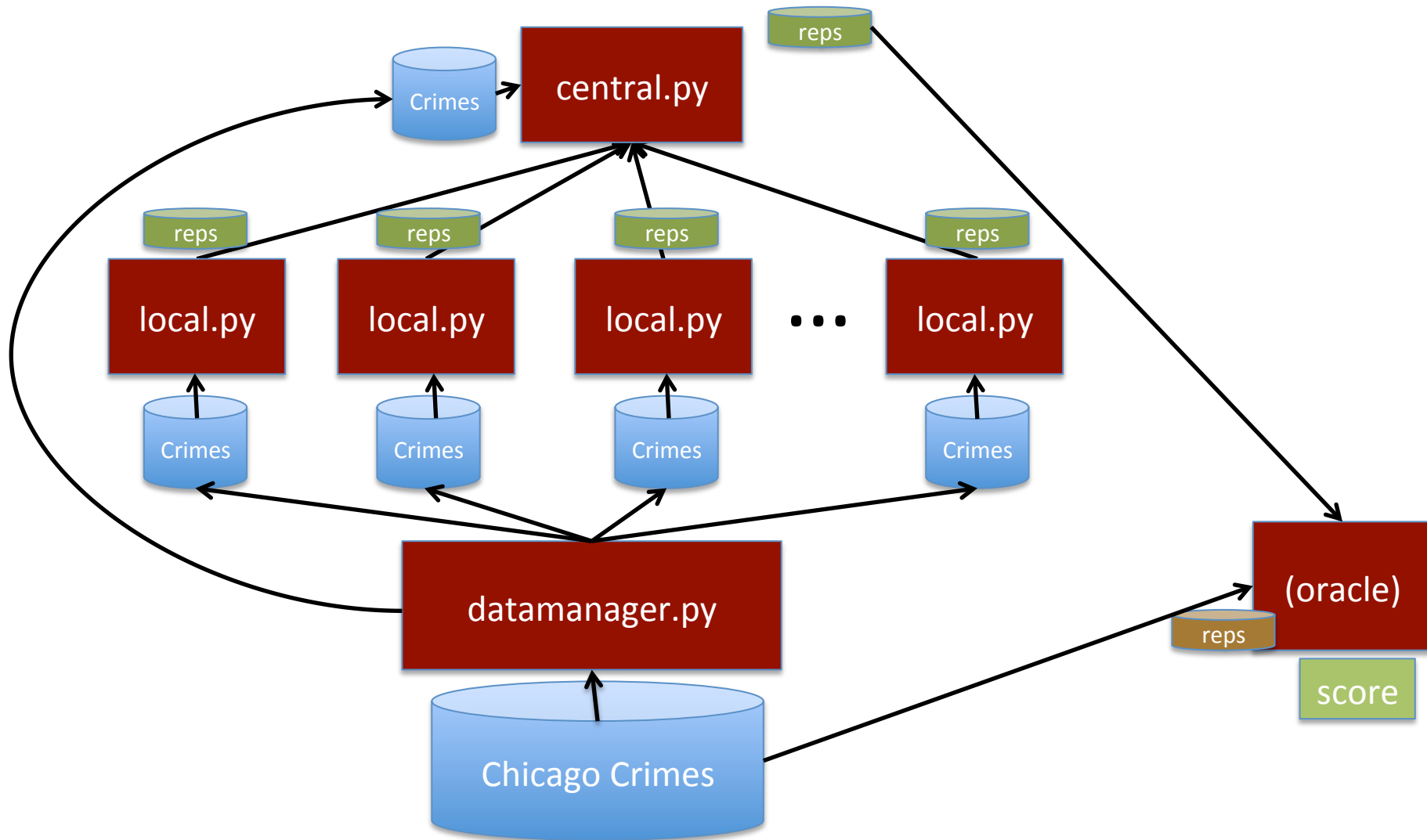
# SIMULATION



# SIMULATION

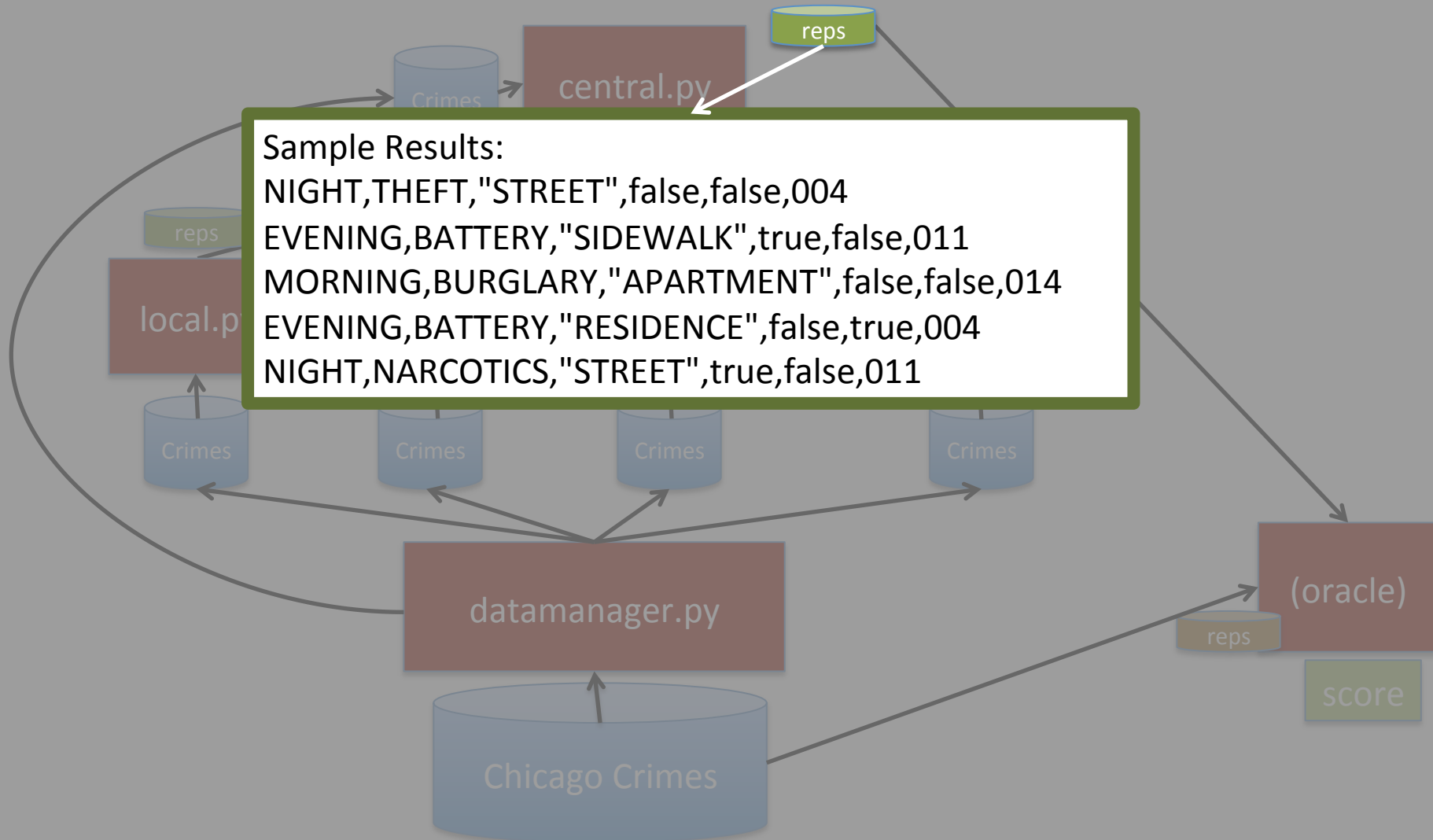


# SIMULATION

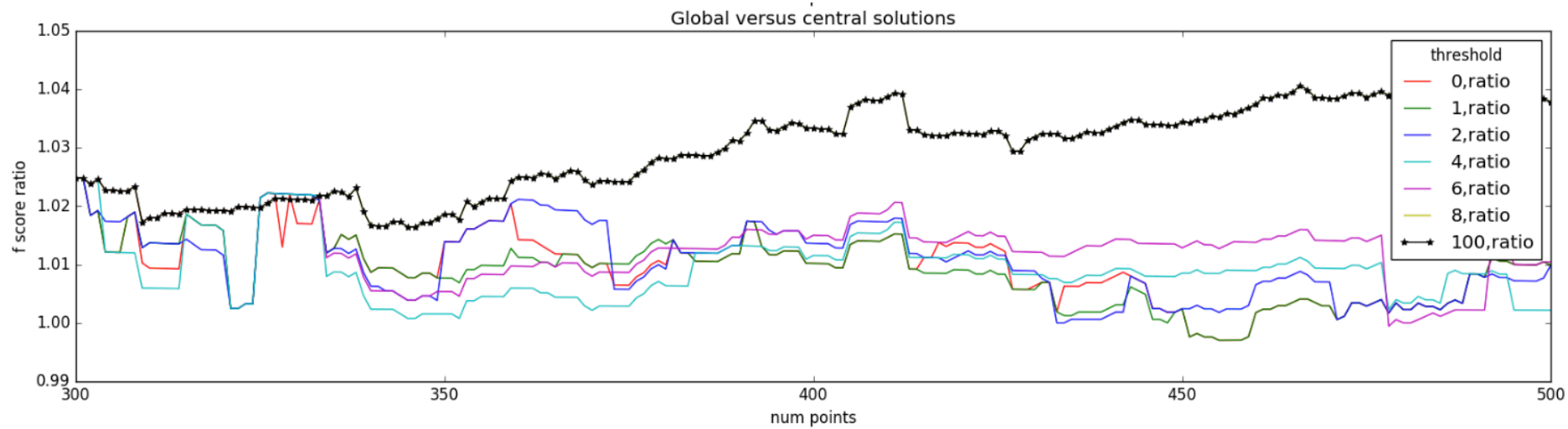




# SIMULATION



# RESULTS



2 local processes, 10 representatives, 300 to 430 points



# NEXT STEPS

- Immediate:
  - Characterize tradeoff
  - Simulate failures
- General:
  - Theoretical bounds
  - Civic applications

