# Towards Human-Enhanced Data Management Systems

**Alessandro Bozzon** - TU Delft
Software and Computer Technology Department
Faculty of Electrical Engineering, Mathematics, and Computer Science (EEMCS)

**Website:** www.alessandrobozzon.com

**E-mail:** a.bozzon@tudelft.nl

**Twitter:** @aleboz

# How can data management systems meet the challenges of next generation knowledge- and data-intensive applications?
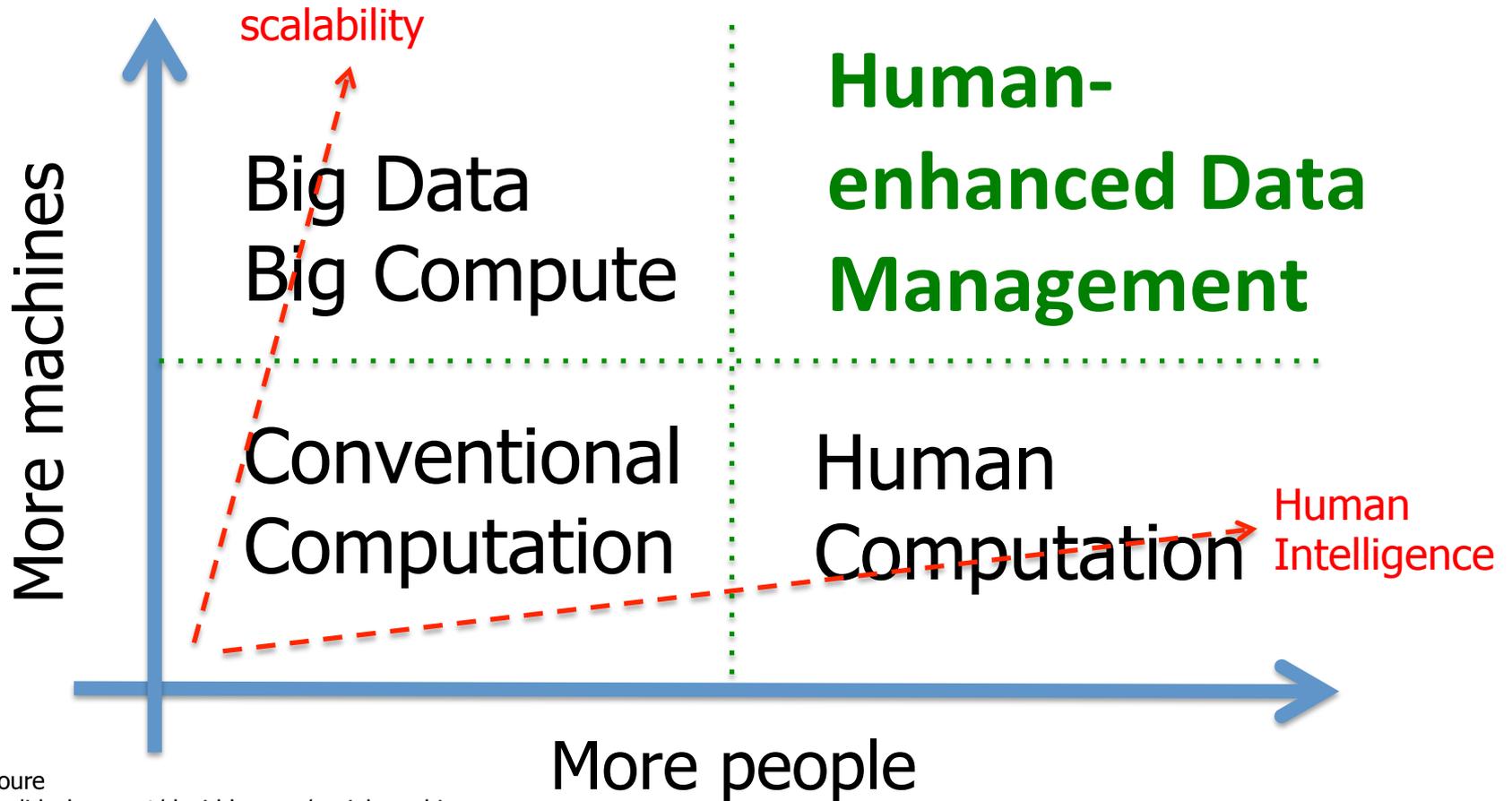
# Big (Ambiguous) Data

Flash Mob Vs. Riot



http://www.telegraph.co.uk/technology/facebook/4542840/Flash-mob-mimicks-T-Mobile-advert-and-closes-train-station.html

# Big (False) Data



Chile Earthquake 2010

# Big Data (Smart) Sense-making

scalability

**Big Data Big Compute**

**Human-enhanced Data Management**

**Conventional Computation**

**Human Computation**

Human Intelligence

More machines

More people

David De Roure
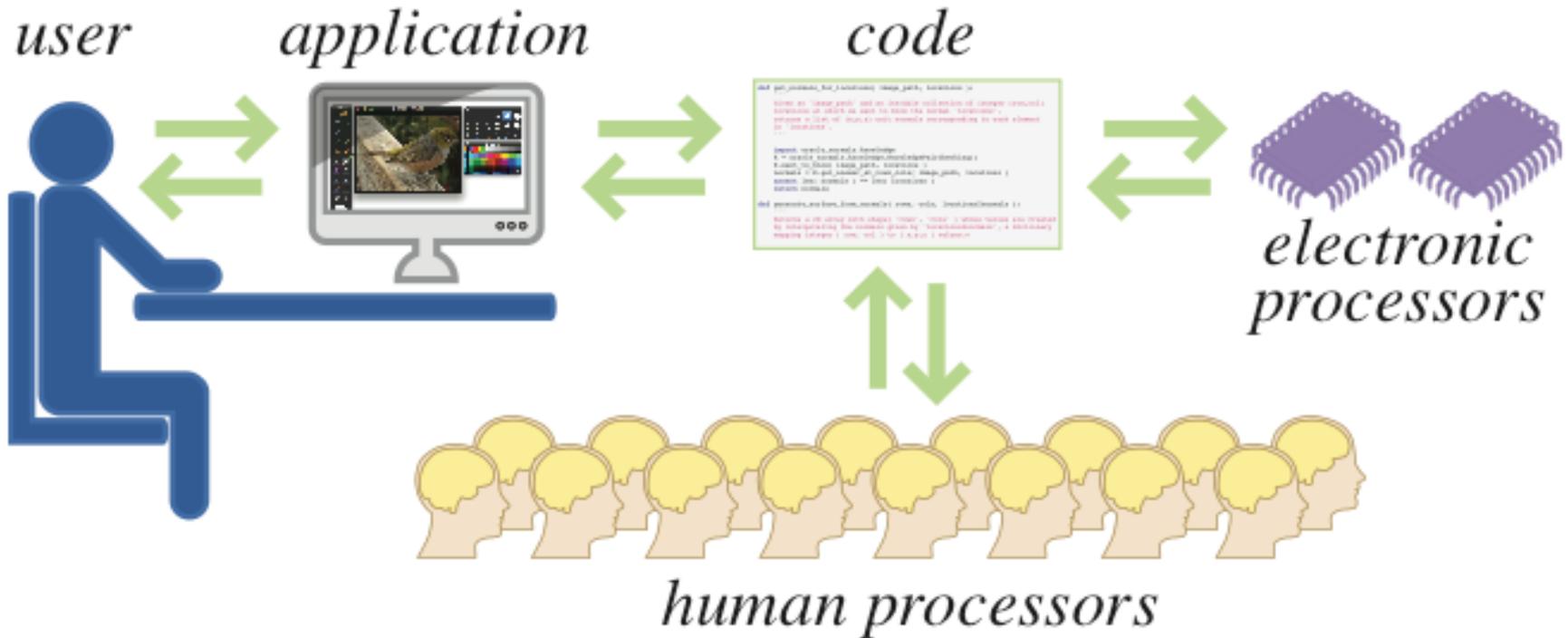http://www.slideshare.net/davidderoure/social-machinesgss

# What do we need?

- A mature theory/practice of automatic <u>and</u> human data management

- A better understanding of the available workforce

- A better understanding of workers engagement mechanisms

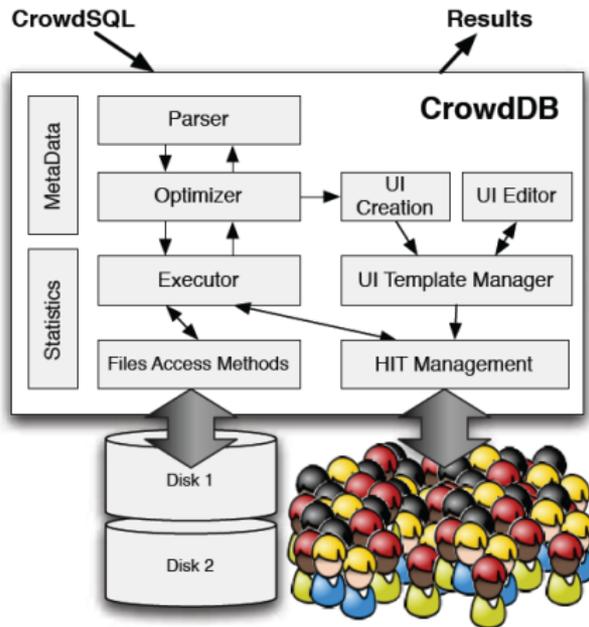# A mature theory/practice of automatic <u>and</u> human data management

# A theory of hybrid computation

## The human co-processing units (HPU)



user   application   code

electronic processors

human processors

- Abstractions?
- Instruction set?
- Design tools?
- Debug?
- Control?

TUDelft

# CrowdDB



- CROWD columns
  - entities known, properties of entities may be unknown
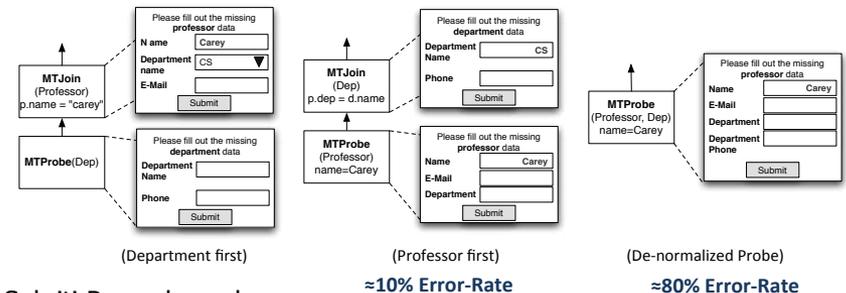
  CREATE TABLE company (
  name STRING PRIMARY KEY,
  hq_address **CROWD** STRING);

- CROWD table
  - entities unknown, crowd-source new entities

  CREATE **CROWD** TABLE department (
      university STRING,
      department STRING,
      phone_no STRING)
      PRIMARY KEY(university, department) );



- GOAL: crowd-source comparisons, missing data
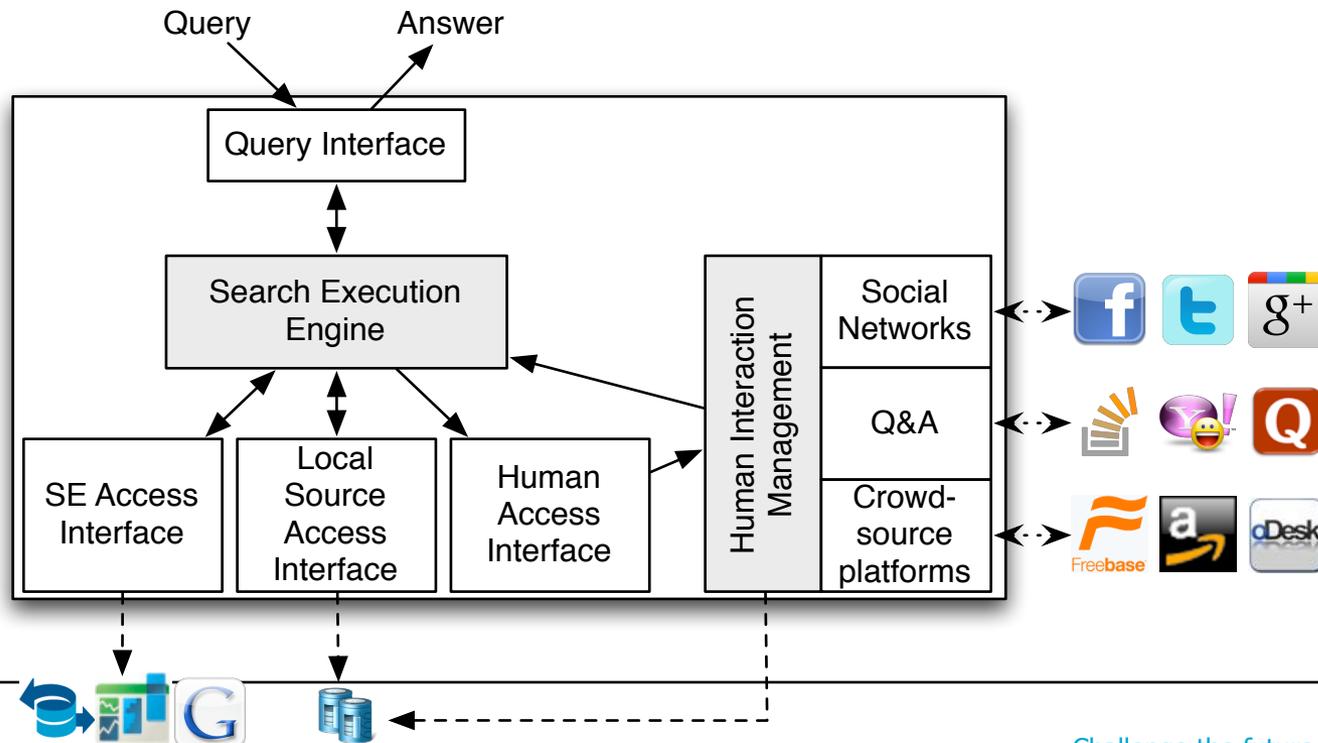  - SQL with extensions to the DML and the Query Language

Michael J. Franklin, Donald Kossmann, Tim Kraska, Sukriti Ramesh, and Reynold Xin. 2011. CrowdDB: answering queries with crowdsourcing. In Proceedings of the 2011 ACM SIGMOD International Conference on Management of data (SIGMOD '11). ACM, New York, NY, USA, 61-72.
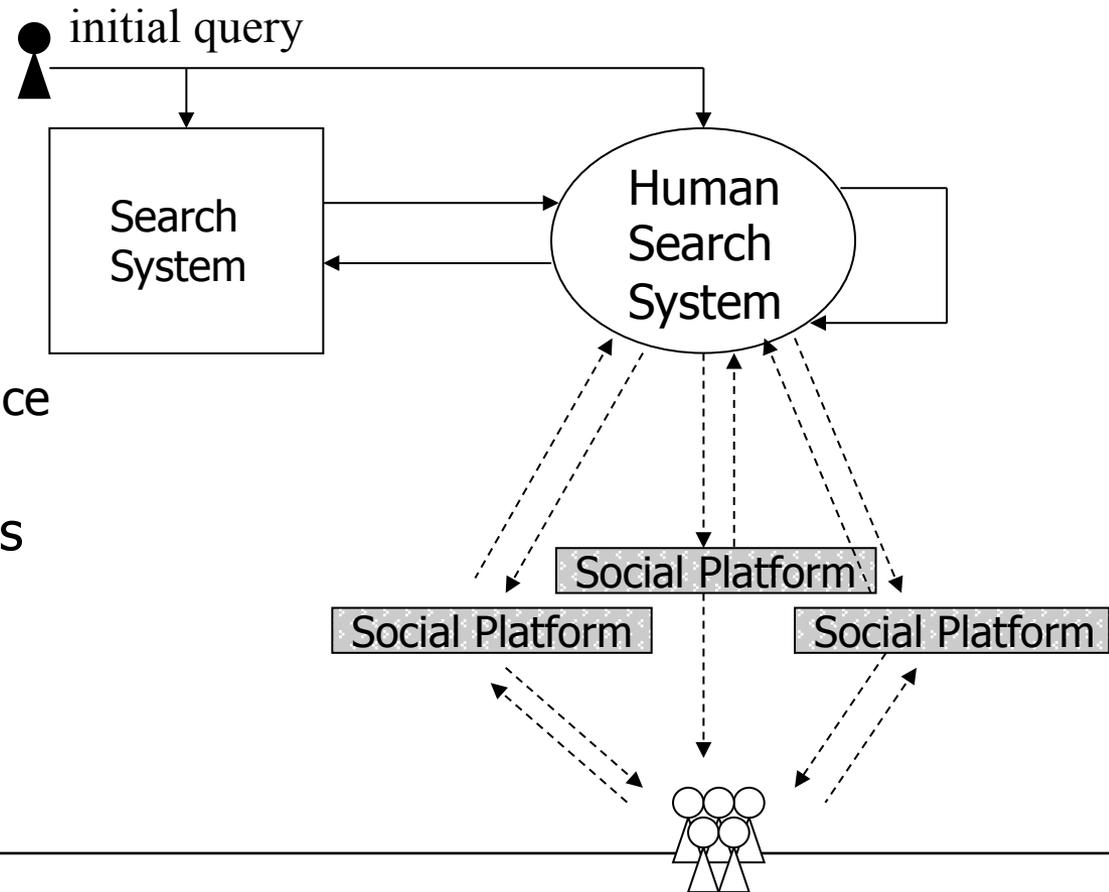
# CrowdSearcher

[Bozzon2012,WWW][Bozzon2013,WWW]

- Multi-platform, reactive, social-network-enabled crowdsourcing
- Our approach: a coordination engine which keeps an overall control on the application deployment and execution

# An example of crowd-based application: crowd-search

- **People do not trust web search completely**
  - Want to get direct feedback from people
  - Expect recommendations, insights, opinions, reassurance

- From search results to friends and experts feedback

initial query

Search System

Human Search System

Social Platform

Social Platform

Social Platform

TUDelft

# Example: Find your next job (exploration)

# Example: Find your job (social invitation)

# Example: Find your job (social invitation)



**Selected data items
can be transferred
to the crowd question**

# Find your job (response submission)



TUDelft

# Crowdsearcher results (in the loop)

# Deployment: search on the social network

- Multi-platform deployment

# Deployment: search on the social network

- Multi-platform deployment

# Deployment: search on the social network

- Multi-platform deployment

# Deployment: search on the social network

- Multi-platform deployment

# Design Process

- A simple task design and deployment process, based on specific data structures
  - created using model-driven transformations
  - driven by the task specification



- **Task Specification:** task operations, objects, and performers
- **Task Planning:** work distribution
- **Control Specification**: task control policies

# Design Dimensions

- Support for several types of task operations
  - Like, Comment, Tag, Classify, Add, Modify, Order, etc.

- Several strategies for
  - **Task splitting**: the input data collection is too complex relative to the cognitive capabilities of users.
  - **Task structuring**: the query is too complex or too critical to be executed in one shot.
  - **Task routing**: a query can be distributed according to the values of some attribute of the collection
  - Output aggregation

- Platform/community assignment
  - a task can be assigned to different communities or social platforms based on its focus

- Control => time, quality, money
  - A reactive execution environment for requirement enforcement and reactive execution

# Reactive Crowdsourcing

- A **conceptual framework** for controlling the execution of crowd-based computations. Based on:
  - Control Marts
  - Active Rules

- Classical forms of controls:
  - Majority control (to close object computations)
  - Quality control (to check that quality constraints are met)
  - Spam detection (to detect / eliminate some performers)
  - Multi-platform adaptation (to change the deployment platform)
  - Social adaptation (to change the community of performers)

- Why Active Rules?

# Auxiliary Structures

- **Object :** tracking object responses
- **Performer:** tracking performer behavior (e.g. spammers)
- **Task:** tracking task status

# Rule Example

```
e: UPDATE FOR μTaskObjectExecution[ClassifiedParty]
c: NEW.ClassifiedParty == 'Republican'
a: SET ObjectControl[oID == NEW.oID].#Eval+= 1
```

# Rule Example

```
e: UPDATE FOR µTaskObjectExecution[ClassifiedParty]
c: NEW.ClassifiedParty == 'Republican'
a: SET ObjectControl[oID == NEW.oID].#Eval+= 1
```

# Rule Example

```
e: UPDATE FOR μTaskObjectExecution[ClassifiedParty]
c: NEW.ClassifiedParty == 'Republican'
a: SET ObjectControl[oID == NEW.oID].#Eval+= 1
```

# CROWD-POWERED SEARCH



- Users ask questions on Twitter
- An hybrid system provide answers
- Workers used for
  - label tweets as "rhetorical question" or not
    - Median 3.02 mins
  - produce responses to question
    - Median 77.4 mins
  - Voting responses
    - Median 82.1 mis

**Median time** =>162.5 minutes

**Cost** => $0.95 per tweet

A Crowd-Powered Socially Embedded Search Engine. Jin-Woo Jeong, Meredith Ringel Morris, Jaime Teevan, Daniel Liebling. ICWSM 2013

# Crowdweaver

- [Kittur et al. 2012]

# Better Understanding of the Available Workforce

# Worker Modeling



- Implicit/Explicit Knowledge
- Interests
- Hard/Soft Skills

- Attitude
- Task-type/domain specific performance
- Availability

- Capacity
- Teams/Groups
- Trust

# Problem

- Ranking the members of a social group according to the level of knowledge that they have about a given topic
- Application: crowd selection (for Crowd Searching or Sourcing)

- Available data
  - User profile
  - behavioral trace that users leave behind them through their social activities

TUDelft

# Finding the right crowd

[Bozzon2013,EDBT]

- Ranking the members of a social group according to the level of knowledge that they have about a given topic
  - Application: crowd selection (for Crowd Searching or Sourcing)

- Available data
  - User profile
  - behavioral trace that users leave behind them through their social activities

# Main Results

- **Profiles** are **less effective** than level-1 resources
  - Resources **produced by others help** in describing each individual's expertise

- **Twitter** is the **most effective** social network for expertise matching – sometimes it outperforms the other social networks
  - Twitter most effective in Computer Engineering, Science, Technology & Games, Sport

- **Facebook** effective in Locations, Sport, Movies & TV, Music

- **Linked-in** never very helpful in locating expertise

# Trustworthiness of Social Information

- Strong correlation between the number of resources and the retrieval performance
- Users DO NOT expose all their interests on social networks

# Better understanding of worker engagement mechanisms

# Critical Mass



**Workers**                    **Work**

# Worker Engagement

## "money, love, or glory"





I ❤ Crowd Sourcing



T. W. Malone, R. Laubacher, and C. Dellarocas. Harnessing Crowds: Mapping the Genome of Collective Intelligence. Working paper no. 2009-001, MIT Center for Collective Intelligence, Feb. 2009.

# Thank You

## **Alessandro Bozzon** - TU Delft

Software and Computer Technology Department

Faculty of Electrical Engineering, Mathematics, and Computer Science (EEMCS)

**Website:** www.alessandrobozzon.com

**E-mail:** a.bozzon@tudelft.nl

**Twitter:** @aleboz

# Rule Programming Best Practice

- We define **three** classes of rules
  - → **Control rules**: modifying the *control* tables;
  - ⇢ **Result rules:** modifying the *dimension* tables (object, performer, task);
  - ⋯⋗ **Execution rules:** modifying the *execution* table, either directly or through re-planning



- Termination must be proven (rule precedence graph has cycles)

# Crowdsearcher Experiment 1

- **Goal:** Test engagement on social networks
- Some 150 users
- Two classes of experiments:
  - Random questions on fixed topics: interests (e.g. restaurants in the vicinity of Politecnico), to famous 2011 songs, or to top-quality EU soccer teams
  - Questions manually submitted by the users
- Different invitation strategies:
  - Random invitation
  - Explicit selection of responders by the asker
- Outcome
  - 175 *like* and *insert* queries
  - 1536 invitations to friends
  - 230 answers
  - 95 questions (~55%) got at least one answer

**TU**Delft

# Manual and Random Questions

# Interest / Rewarding Factor

- Manually written and assigned questions are consistently more responded in time

# Query Type

- Engagement depends on the difficulty of the task
- Like vs. Add tasks:

TUDelft

# Experiments: Distribution of answers/invitation

TUDelft

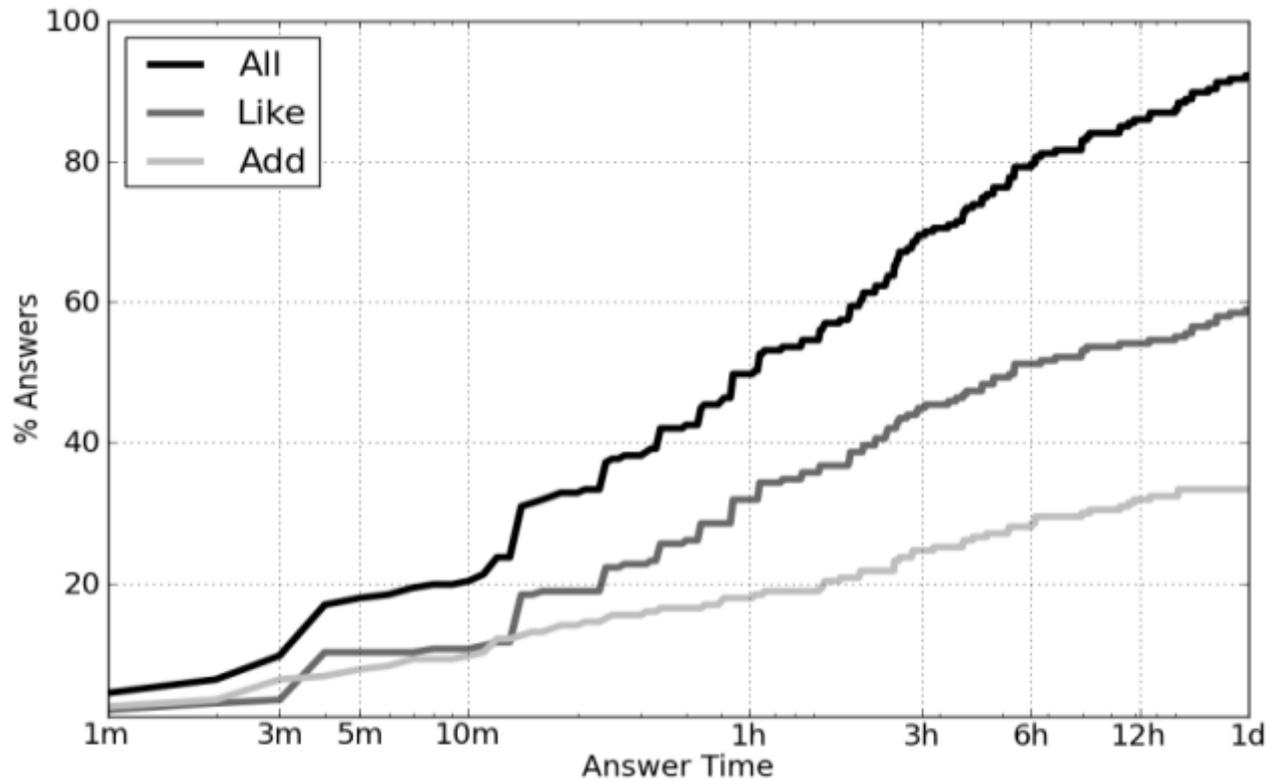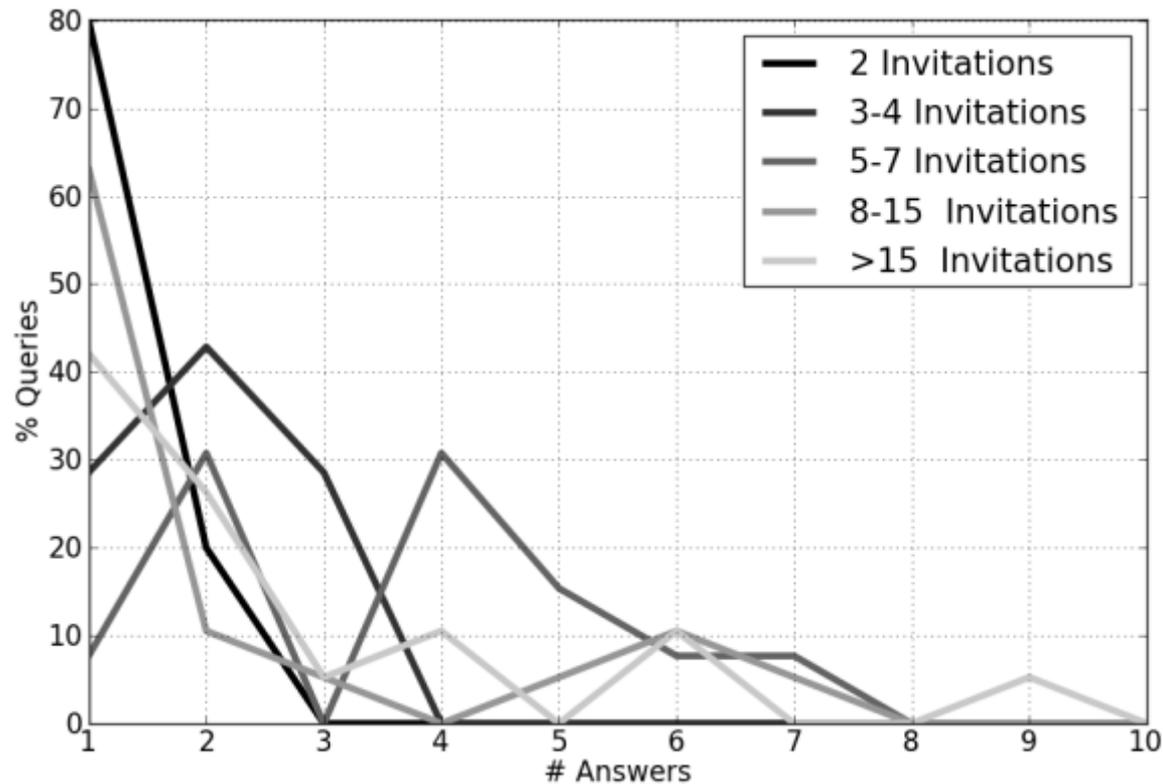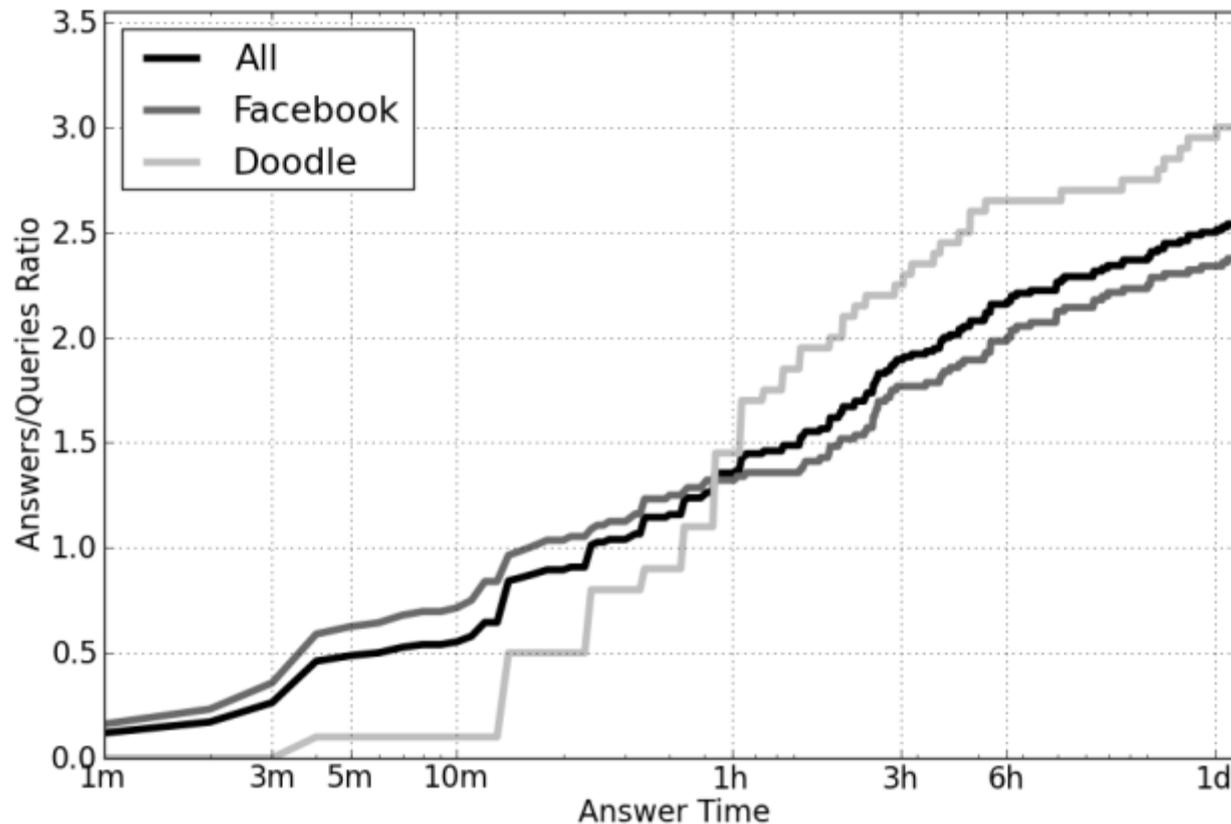Sometimes: more answers than invitations (limited cases)

# Comparison of Execution Platforms

- Facebook vs. Doodle

# Posting Time

- Facebook vs. Doodle

# Crowdsearcher Experiment 2

- **GOAL**: demonstrate the *flexibility* and *expressive power* of reactive crowdsourcing

- 3 experiments, focused on Italian politicians
  - **Parties**: Human Computation → affiliation classification
  - **Law**: Game With a Purpose → guess the convicted politician
  - **Order**: Pure Game → hot or not
- 1 week (November 2012)
- 284 distinct performers
  - Recruited through public mailing lists and social networks announcements
- 3500 Micro Tasks

# Politician Affiliation

- Given the picture and name of a politician, specify his/her political affiliation
  - No time limit
  - Performers are encouraged to look up online

- 2 set of rules
  - Majority Evaluation
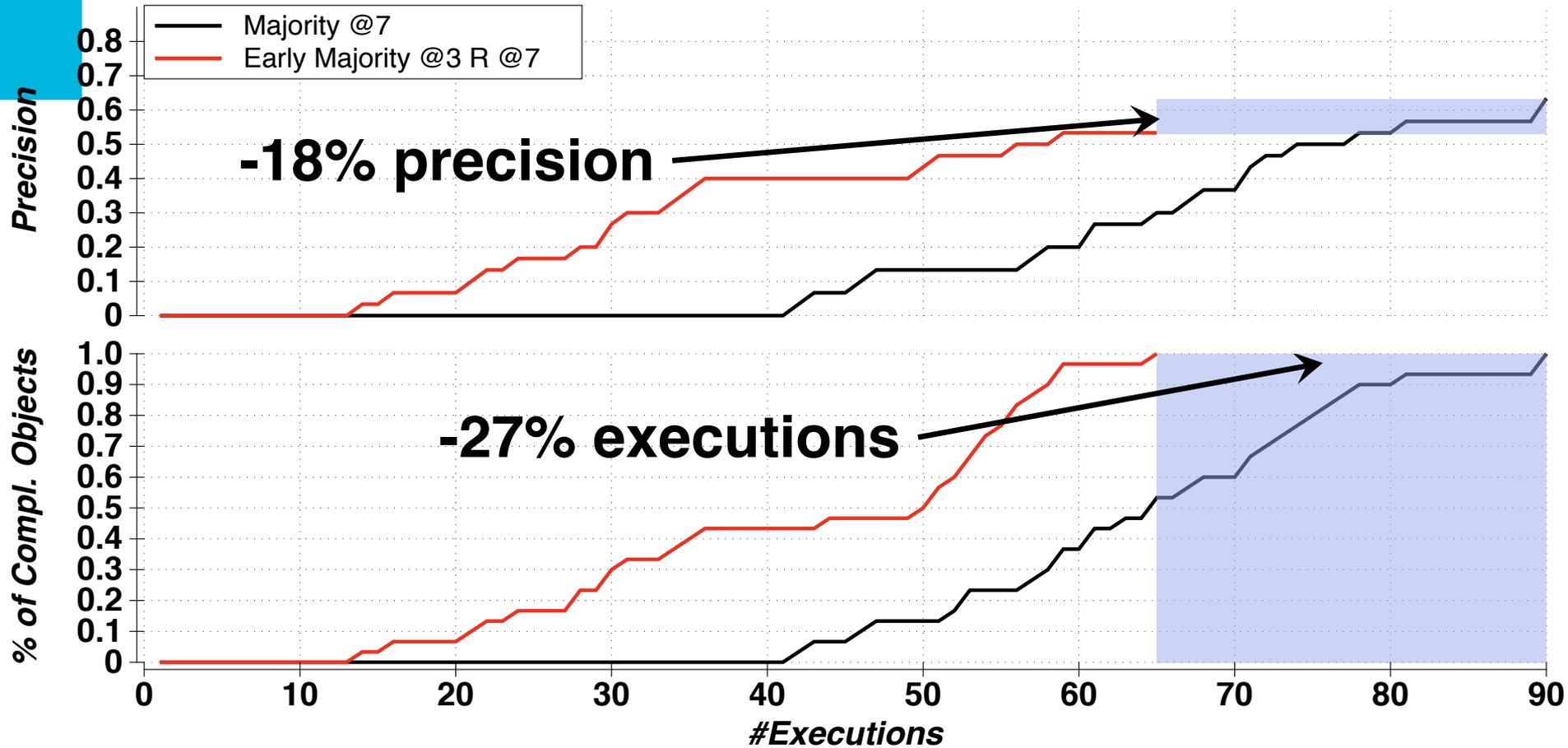  - Spammer Detection



TUDelft

# Results – Majority Evaluation$_{1/3}$



30 object; object redundancy = 9;
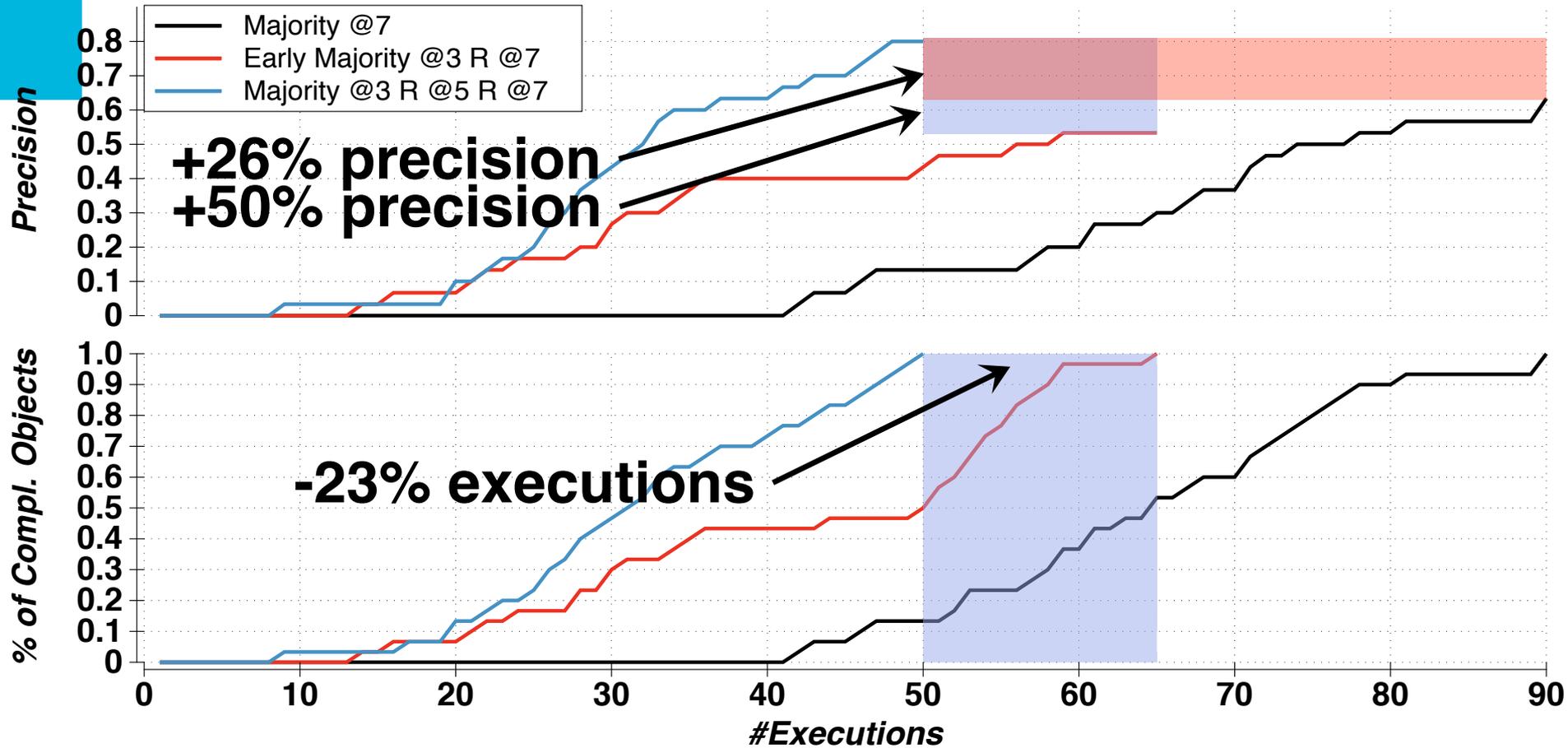Final object classification as simple majority after 7 evaluations
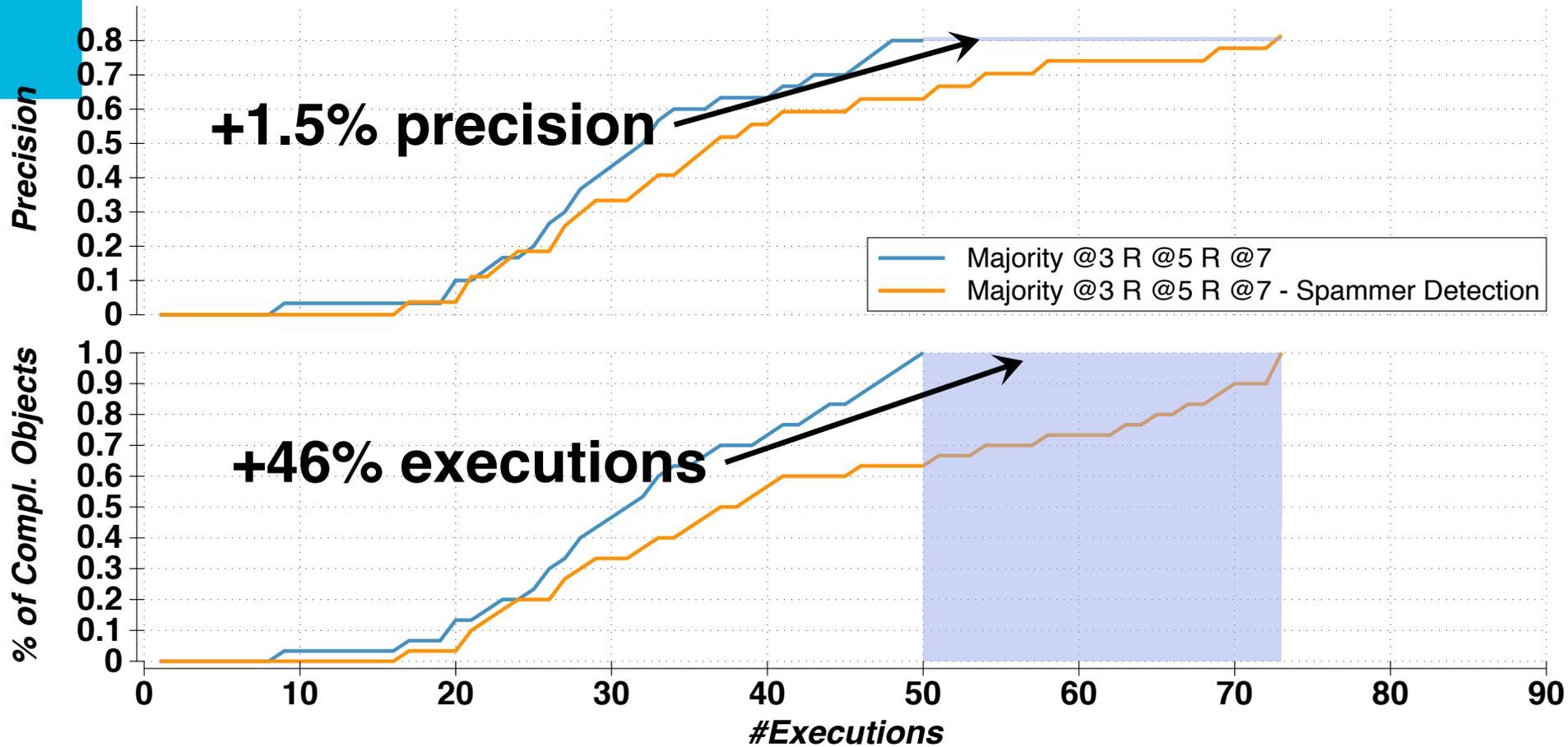
# Results - Majority Evaluation $_{2/3}$



Final object classification as **total** majority after 3 evaluations
Otherwise, re-plan of 4 additional evaluations. Then simple majority at 7

# Results - Majority Evaluation_3/3



Final object classification as **total** majority after 3 evaluations
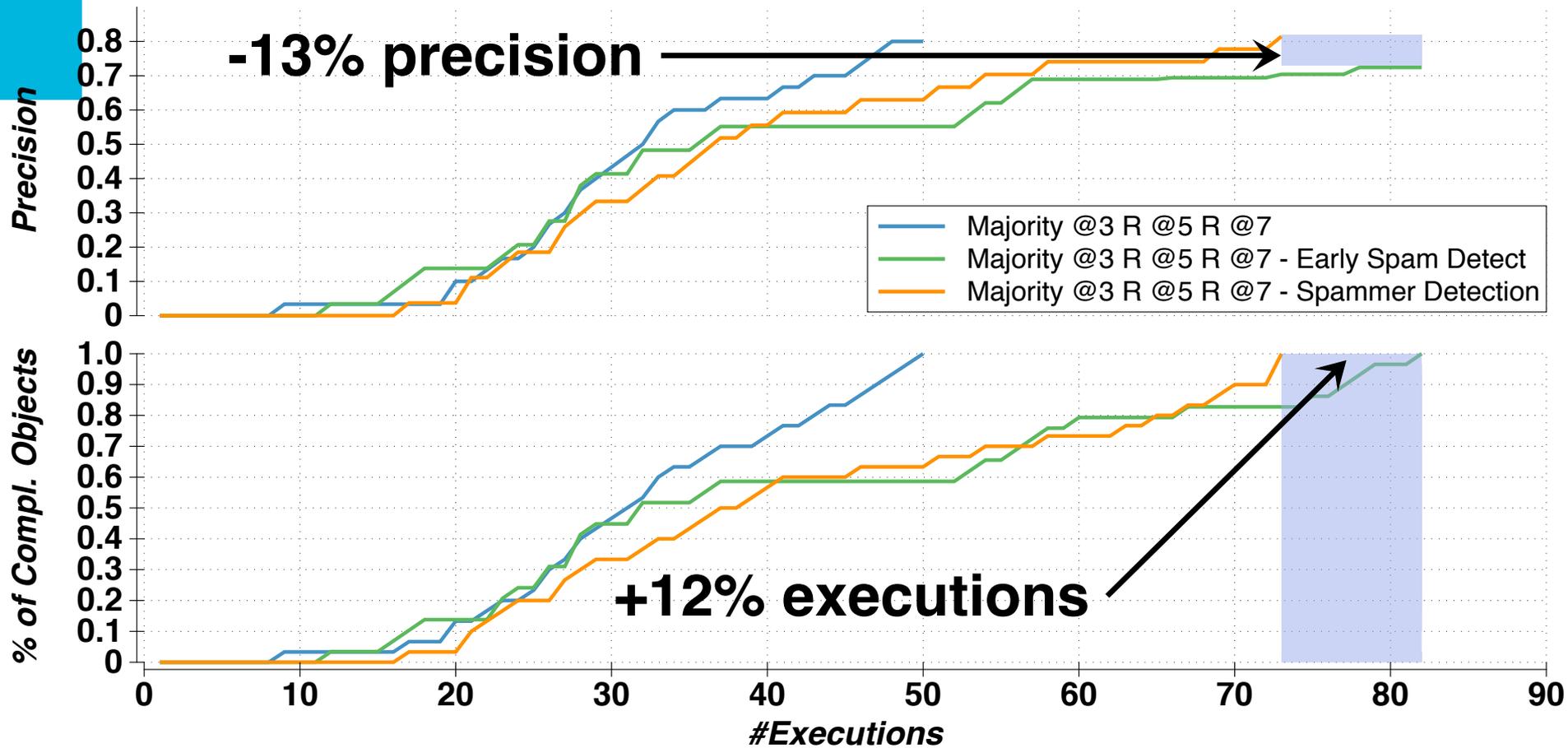Otherwise, **simple** majority at 5 or at 7 (with replan)

# Results – Spammer Detection _1/2_



New rule for spammer detection **without ground truth**
Performer correctness on **final** majority. Spammer if > 50% wrong classifications

# Results – Spammer Detection $_{1/2}$



**-13% precision**

**+12% executions**

**Precision**

**% of Compl. Objects**

**#Executions**

Majority @3 R @5 R @7
Majority @3 R @5 R @7 - Early Spam Detect
Majority @3 R @5 R @7 - Spammer Detection

New rule for spammer detection **without ground truth**

Performer correctness on **current** majority. Spammer if > 50% wrong
classifications