# Causal Discovery from Mass Cytometry Data

*Presenters: Ioannis Tsamardinos and Sofia Triantafillou*

Institute of Computer Science, Foundation for Research and Technology, Hellas
Computer Science Department, University of Crete
in collaboration with Computational Medicine Unit, Karolinska Institutet

# The Measuring Technology
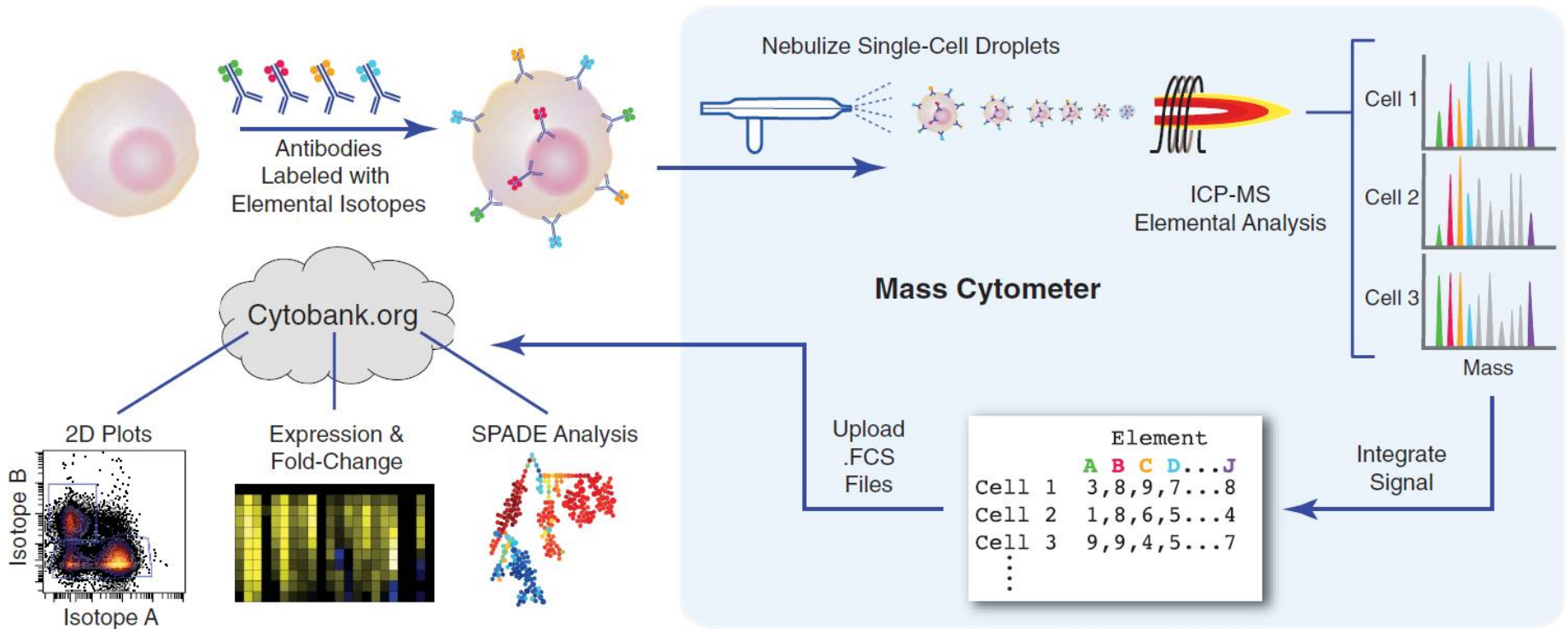
# Mass Cytometry

- **Single cells** measurements

- **Sample sizes in the millions**, minimal cost

- Public data available

- Up to ~30 proteins measured at a time

- Applications
    1. Cell counting
    2. Cell sorting (gating)
    3. Identifying signaling responses
    4. Drug screening
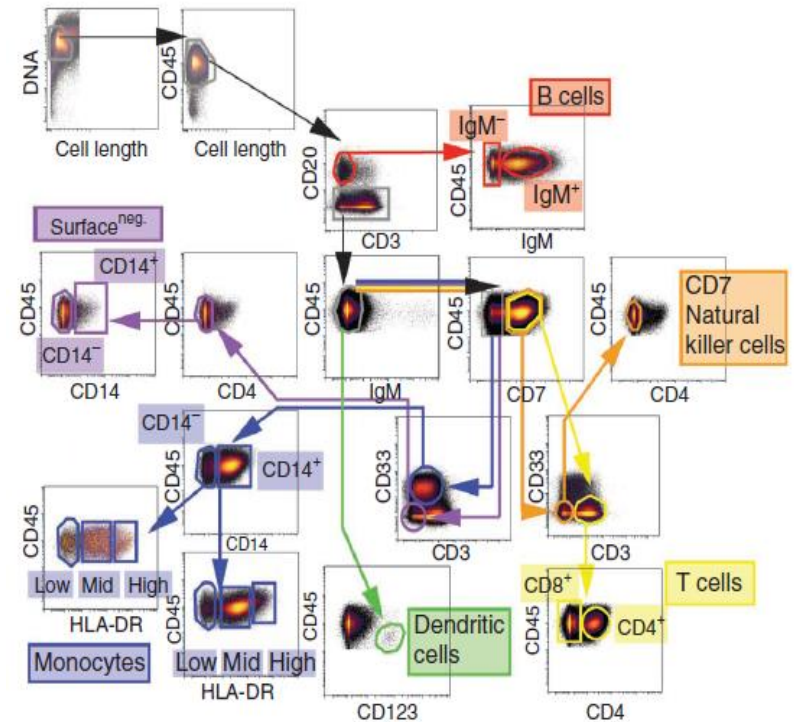    5. De novo, personalized pathway / causal discovery (?)

# Mass Cytometry



[Image by Bendall et al., Science 2011]
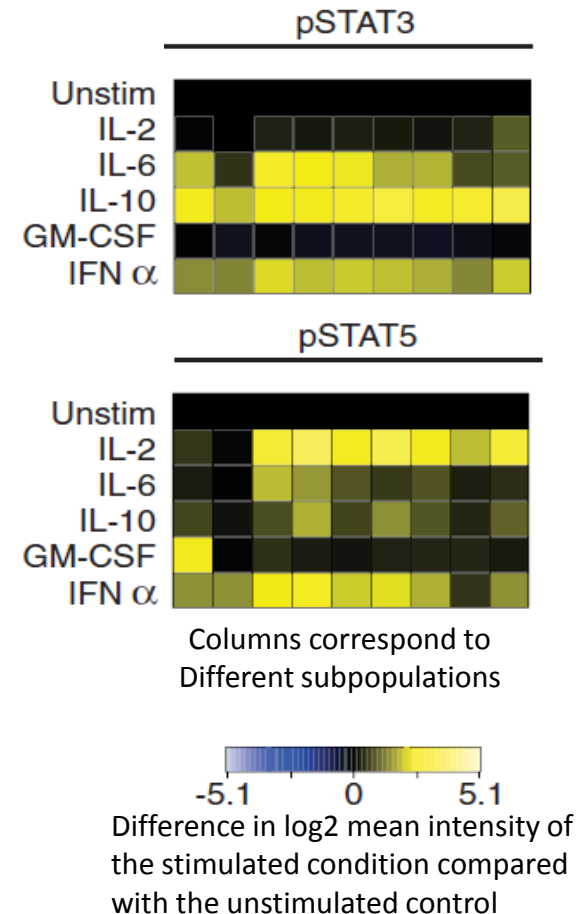
# Cell Sorting (Gating)

- **Immune system cells** can be distinguished based on specific surface markers.

- Process resembles a decision tree



[Image by Bodenmiller et al., Nat. Biotech. 2012]
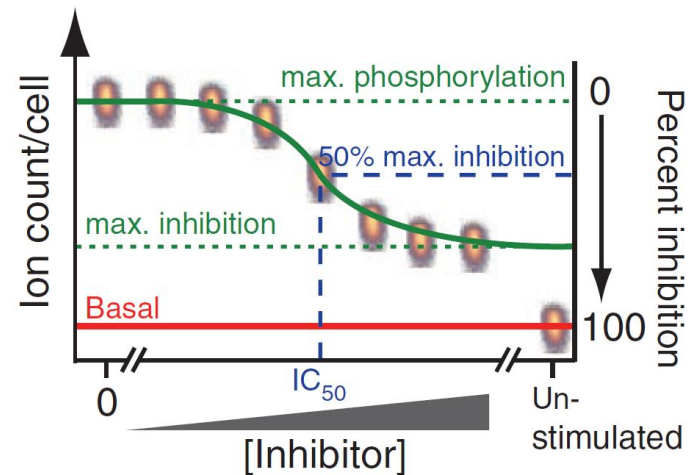
# Identifying Signaling Responses

- Immune responses are triggered by specific **activators**

- Signaling responses are sub-population specific.

- Mass cytometry for identifying signaling effects:

  1. Functional proteins (non-surface) are also marked (e.g., pSTAT3 and pSTAT5)

  2. Activators are applied to stimulate a response to disease

  3. Cells are sorted by sub-population

  4. Changes in protein abundance/phosphorylation in each subpopulation are quantified



Columns correspond to Different subpopulations

Difference in log2 mean intensity of the stimulated condition compared with the unstimulated control

[Image by Bendall et al., Science 2011]

# Drug Screening

- Unwanted signaling responses should be suppressed for disease treatment
- Mass cytometry for drug screening
  1. **After stimulation**, cells are treated with potential drugs (**inhibitors**)
  2. Cells are sorted by sub-population
  3. Dose-response curves are identified
     - Per activator
     - Per sub-population
     - Per inhibitor

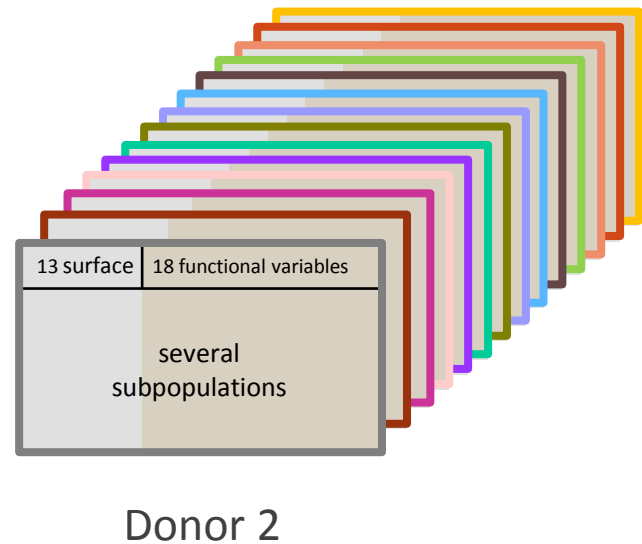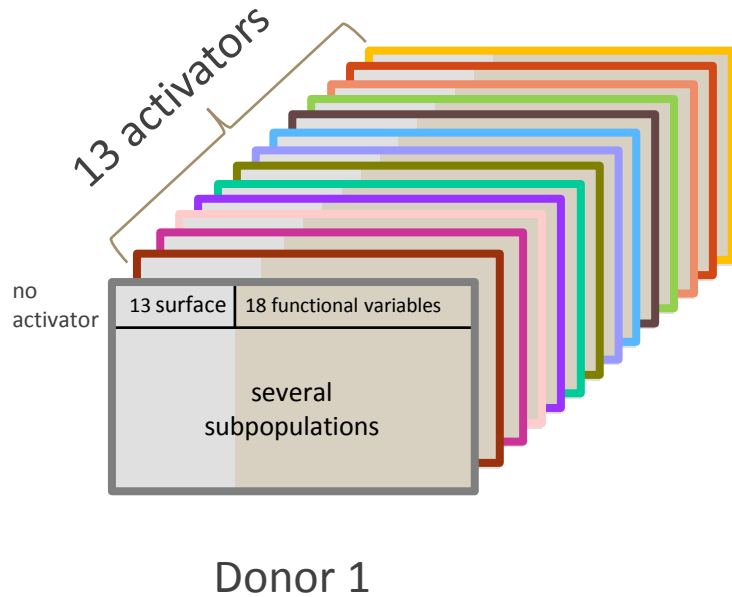[Image by Bodenmiller et al., Nat. Biotech. 2012]

# The Public Data

# Bendall Data
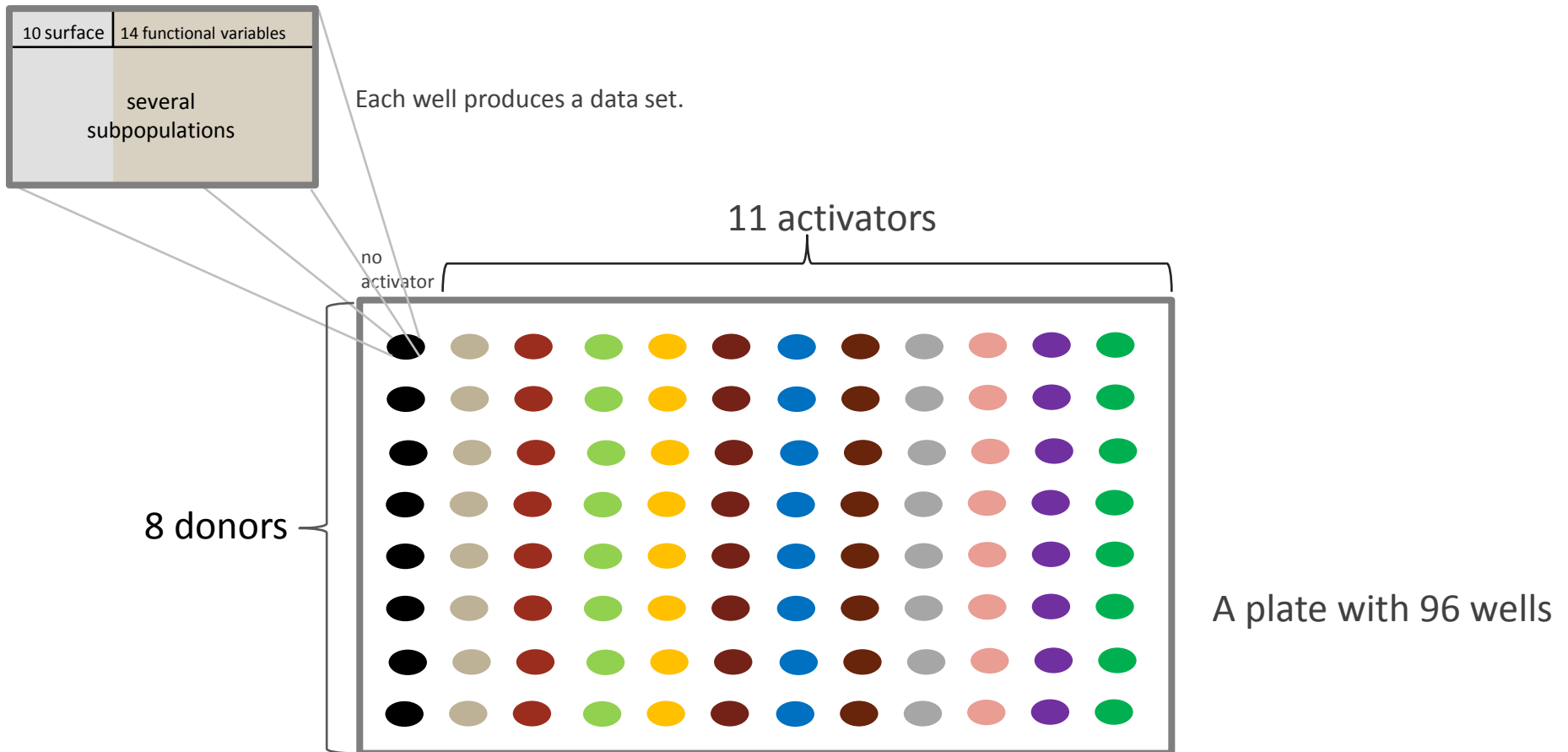


Donor 1

Donor 2

[**Single-Cell Mass Cytometry of Differential Immune and Drug Responses Across a Human Hematopoietic Continuum, Bendall et al.,** *Science* **332**, 687 (2011)]

# Bodenmiller Data: Time Course

10 surface | 14 functional variables

several subpopulations

Each well produces a data set.

11 activators

no activator

0 min
1 min
5 min
15 min
30 min
60 min
120 min
240 min

A plate with 96 wells

[**Multiplexed mass cytometry profiling of cellular states perturbed by small-molecule regulators, Bodenmiller et al.,** *Nature Biotechnology* **30**, 9 (2012) ]

# Bodenmiller Data: 8 donors

10 surface | 14 functional variables

several subpopulations

Each well produces a data set.

11 activators

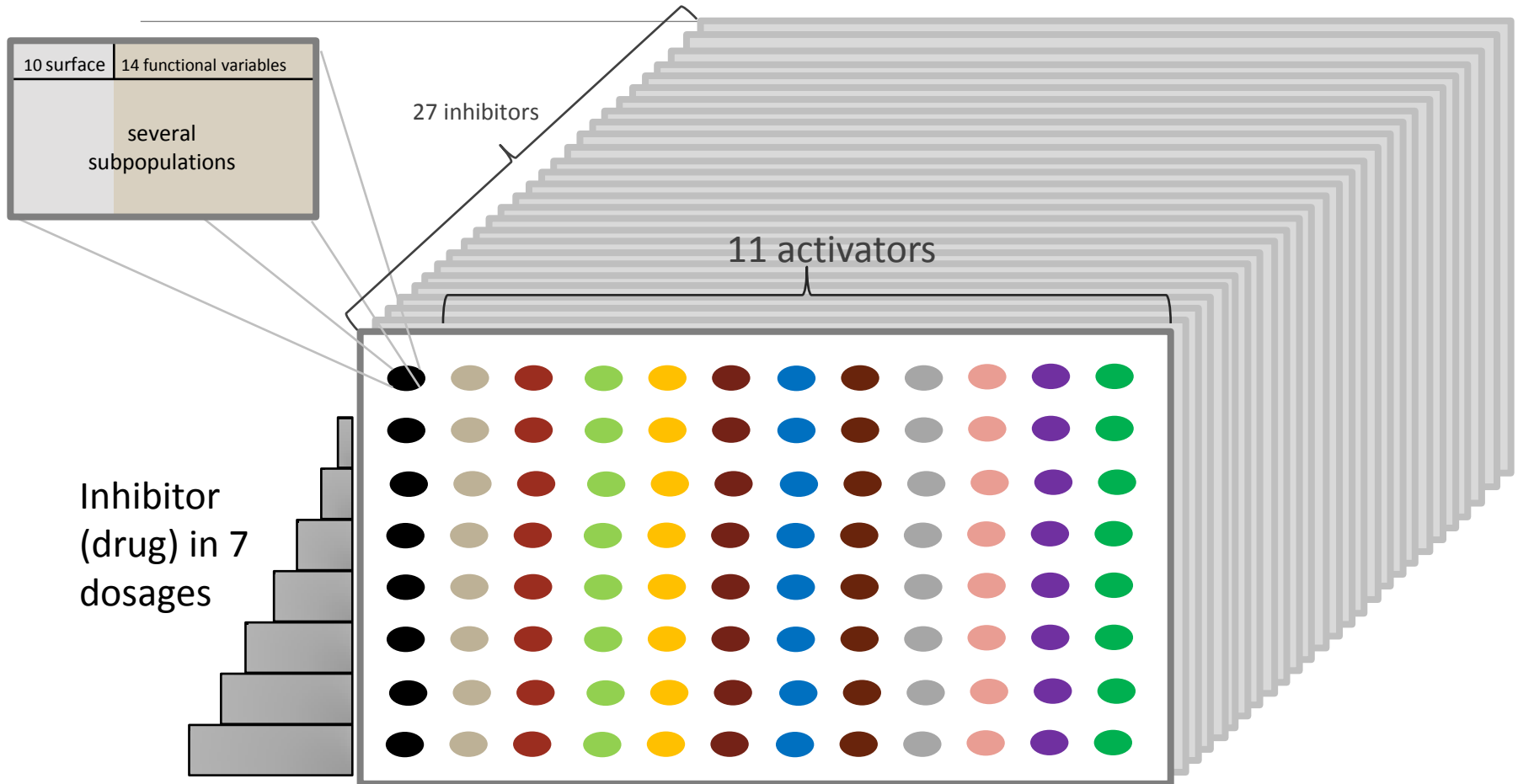no activator

8 donors

A plate with 96 wells

[**Multiplexed mass cytometry profiling of cellular states perturbed by small-molecule regulators, Bodenmiller et al.,** *Nature Biotechnology* **30**, 9 (2012) ]

# Bodenmiller Data: Inhibitors



10 surface | 14 functional variables

several subpopulations

27 inhibitors

11 activators

Inhibitor (drug) in 7 dosages

[**Multiplexed mass cytometry profiling of cellular states perturbed by small-molecule regulators, Bodenmiller et al.,** *Nature Biotechnology* **30**, 9 (2012) ]
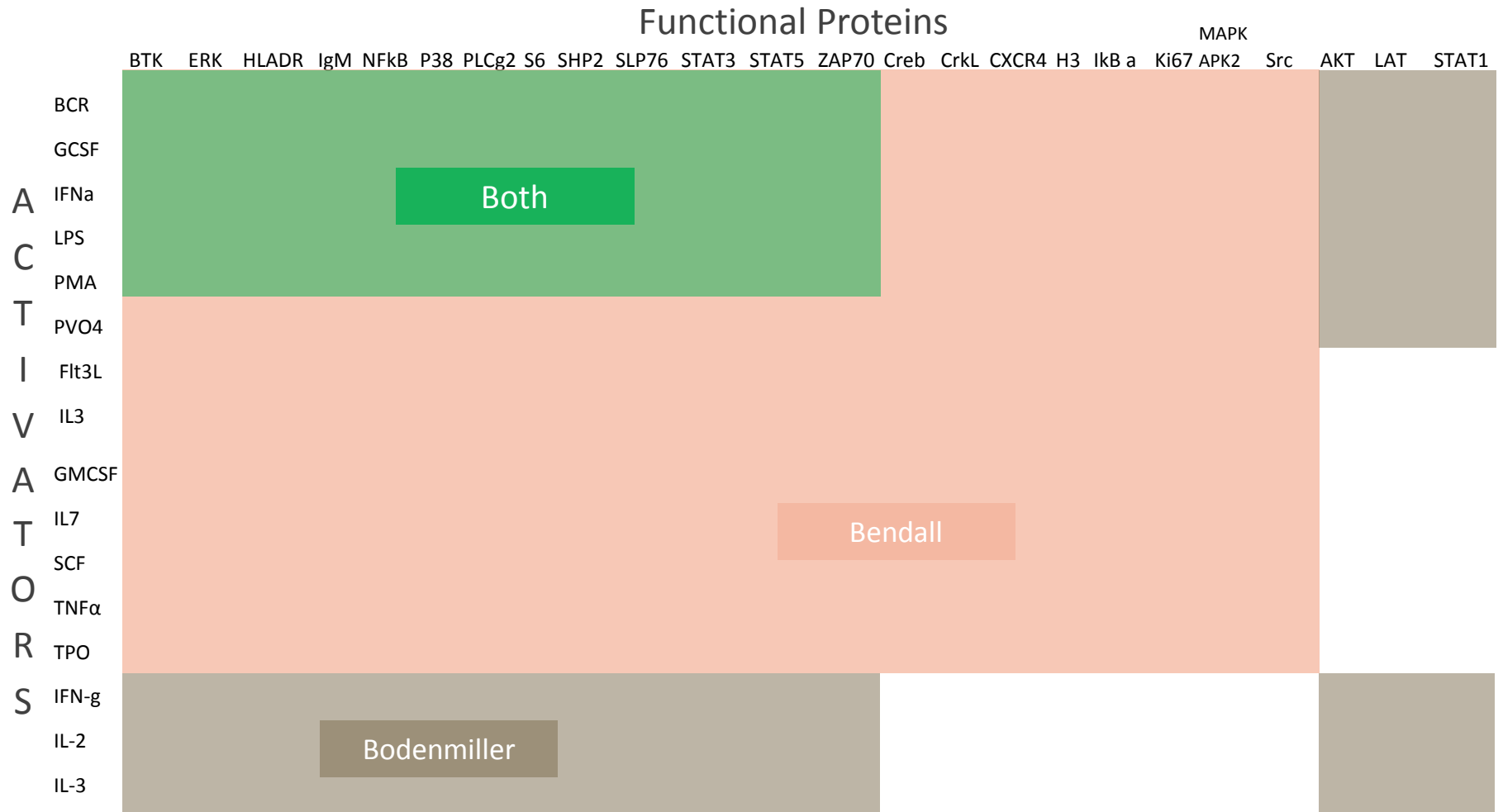
# Data summary

| | Bodenmiller data | | | Bendall data |
|---|---|---|---|---|
| | Inhibitor data | 8donor data | Time course data | |
| Activators | 🟩 | 🟩 | 🟩 | 🟩 |
| Time | | | 🟩 | |
| Donors | | 🟩 | | 🟩 |
| Inhibitors | 🟩 | | | |
| Subpopulations | 🟩 | 🟩 | 🟩 | 🟩 |
| Proteins | 🟩 | 🟩 | 🟩 | 🟩 |

Collection of datasets with :
All activators
1 time point (30')
1 donor
All Inhibitors
All Subpopulations
All 10+14 markers measured

# Data Summary

# Causal Discovery in Mass Cytometry



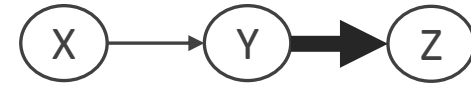Image courtesy of Dr. Brad Marsh

A typical day in the cell

- Feedback loops
- Latent variables
- Non-linear relations
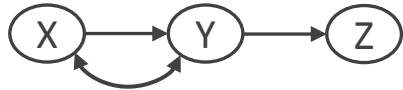- Unfaithfulness

# A Basic Approach

# Local Causal Discovery

# Issue #1:
# Signaling is Sub-Population Specific

- Gate data
  - Data were gated by the initial researchers in Cytobank.org

- Analyze sub-populations independently

- Gated sub-populations differ between Bodenmiller and Bendall
  - cd4+, cd8+, nk sub-populations in common.

| Bodenmiller | | Bendall | | |
|---|---|---|---|---|
| cd14+hladr-, | cd14-surf- | Pre-B II | **Mature CD4+ T** | MPP |
| cd14+hladrhigh | **cd4+** | Mature CD38lo B | **Naive CD4+ T** | HSC |
| cd14+hladrmid | **cd8+** | Pre-B I | CMP | Megakaryocyte |
| cd14+surf- | dendritic | Mature CD38mid B | **Naive CD8+ T** | Erythroblast |
| cd14-hladr- | igm+ | Immature B | **Mature CD8+ T** | Platelet |
| cd14-hladrhigh | igm- | Plasma cell | CD11b- Monocyte | MEP |
| cd14-hladrmid | **nk** | **nk** | CD11bmid Monocyte | Plasmacytoid DC |
| | | Myelocyte | CD11bhi Monocyte | GMP |

# Issue #2:Dormant Relations

- Relations may appear only during signaling
  - Pool together unstimulated and stimulated data

- Different parts of the pathway maybe activated by different activators
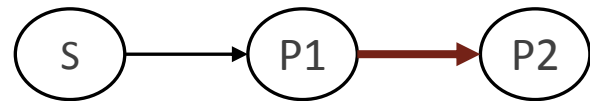  - Analyze data from different activators independently

# Issue #3:
# Testing Independence

- Check (in)dependencies:
    1. $Dep(X, Y | \mathbf{Z})$
    2. $Ind(X, Y | \mathbf{Z})$

S → P1 → P2

- Choosing a test of conditional independence
    ◦ One binary, two continuous variables
    ◦ Relations typically non-linear
    ◦ Options:
        1. Discretization  BUT: does not preserve conditional independencies
        2. Rejected but promising candidates:
            1. Maximal Information Coefficients (Reshef et al., Science 334, 2011)
            2. Kernel-based Conditional Independence test (Zhang et al., UAI 2011)
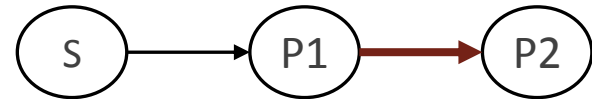        3. Fisher z-test of independence + logistic regression
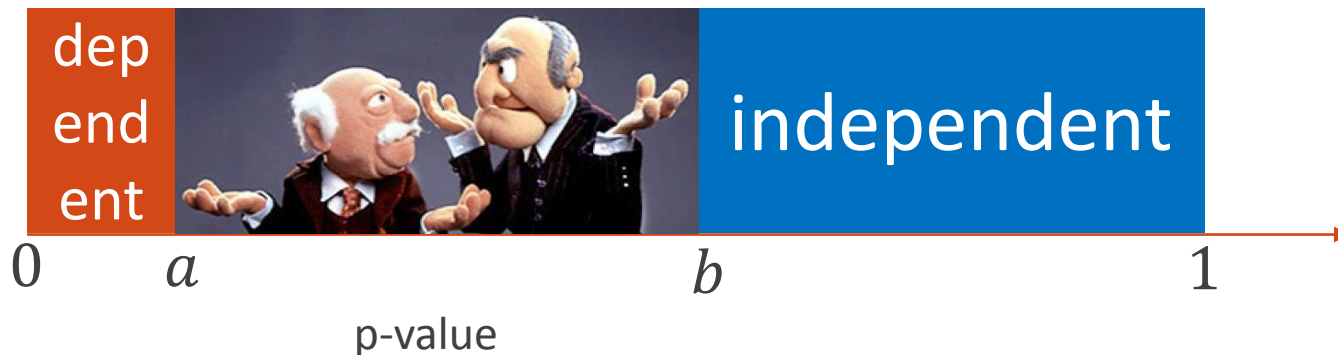
# Issue #4
# Make Reliable Predictions

- Check **ALL** (in)dependencies:
  1. $Dep(S, P1)$
  2. $Dep(S, P2)$
  3. $Dep(P1, P2)$
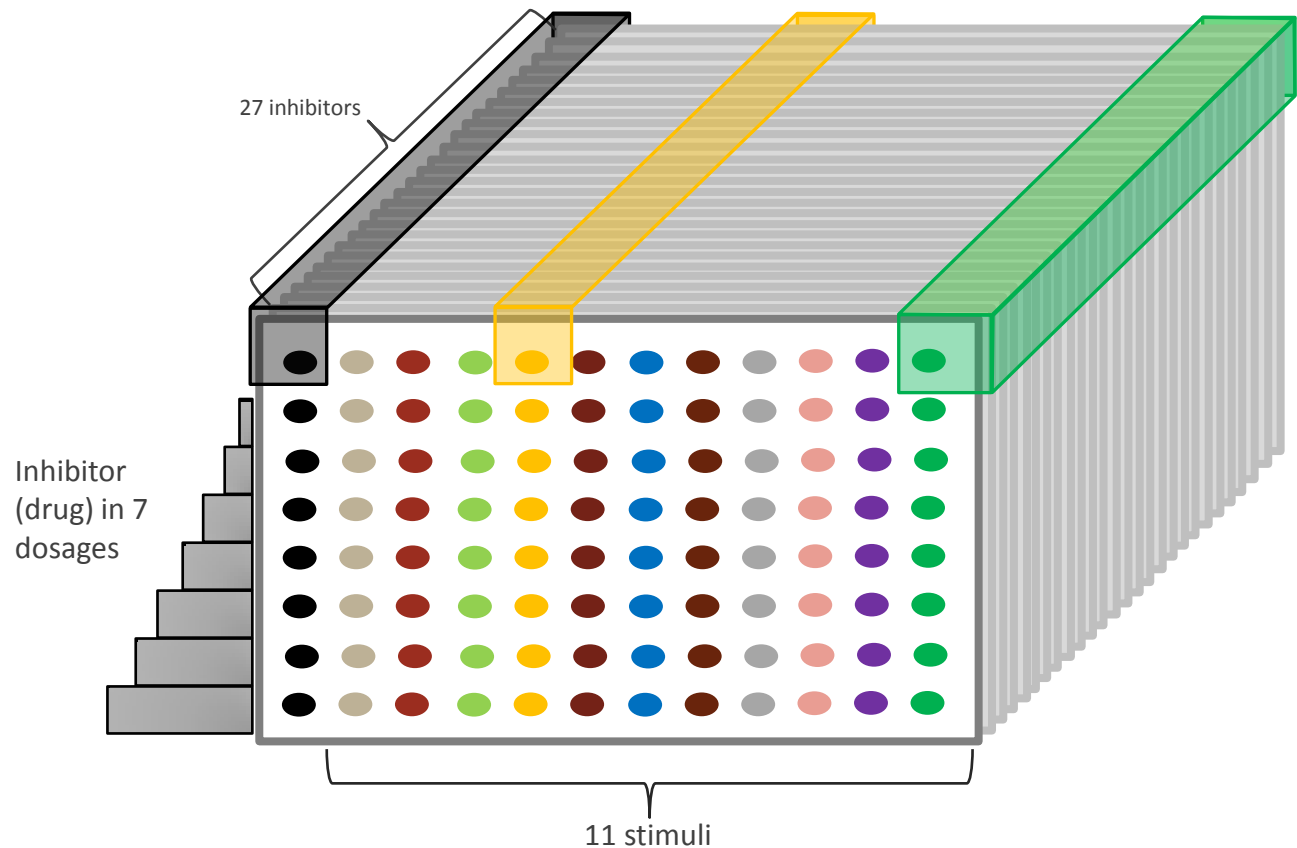  4. $Ind(S, P2|P1)$
  5. $Dep(S, P1|P2)$
  6. $Dep(P1, P2|S)$



- Two thresholds, $a$ =0.05 for dependence, $b$ =0.15 for independence



dep end ent

independent

0    $a$              $b$              1

p-value

# Issue #5:
# Identify "Outlier" Experiments

- Inhibitor data for "zero" dosage and 8 donor data should represent the same joint distribution

- Do they?



27 inhibitors

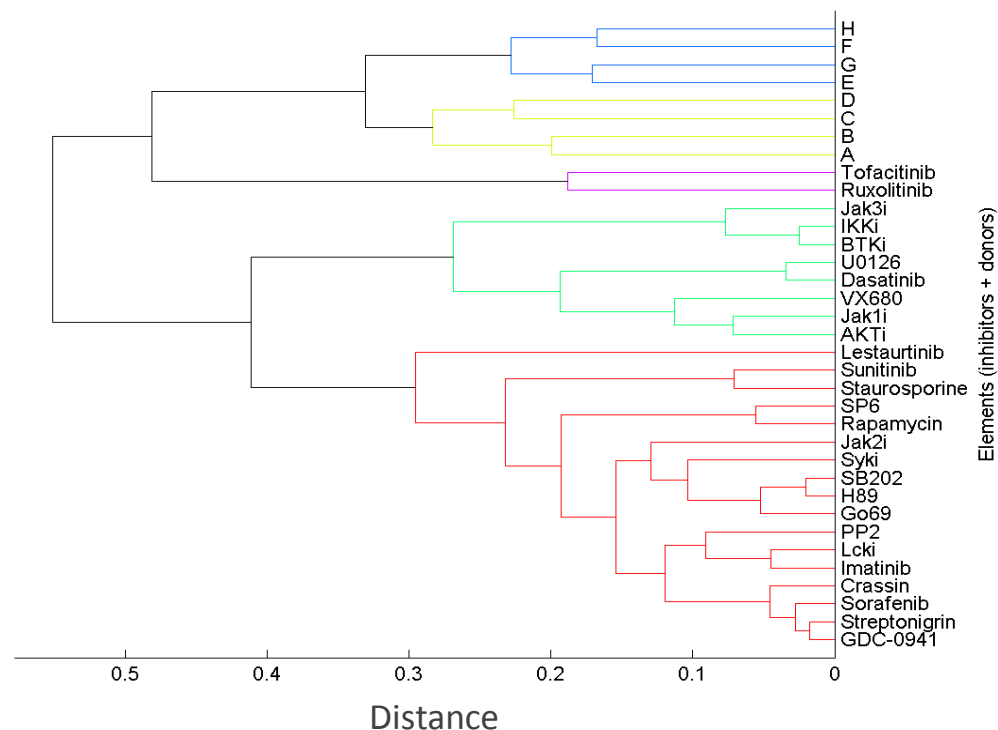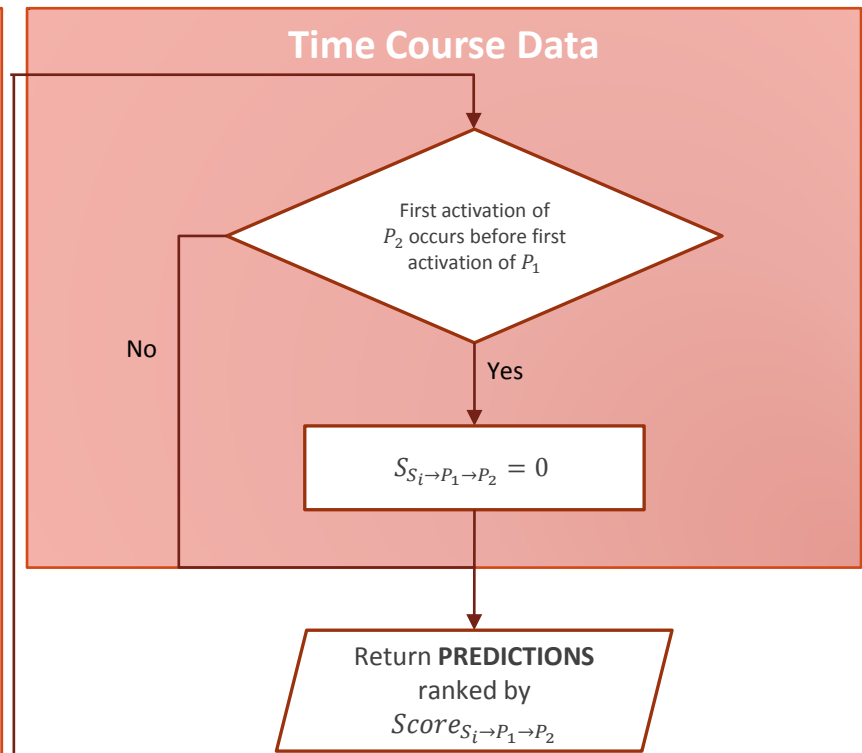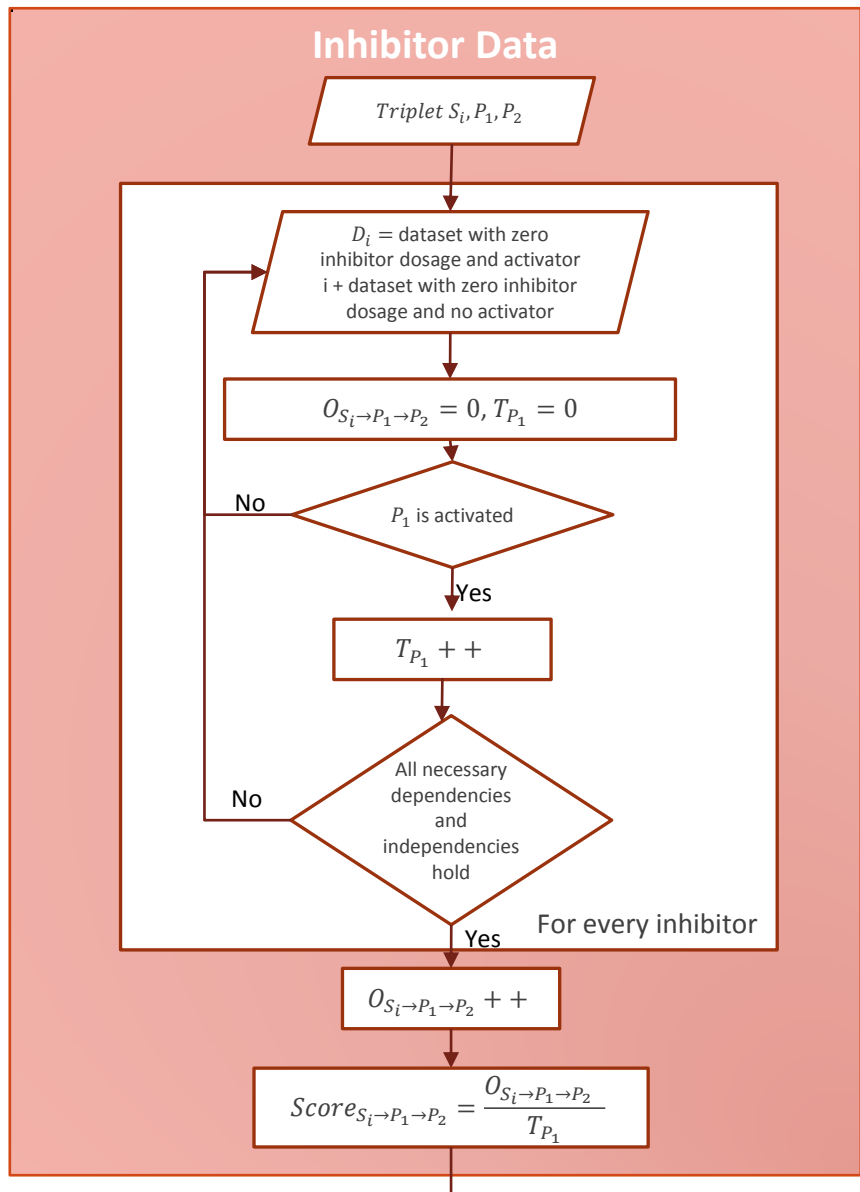Inhibitor (drug) in 7 dosages

11 stimuli

# Issue #5:
# Identify "Outlier" Experiments

- Inhibitor data for "zero" dosage and 8 donor data should represent the same joint distribution

- Do they?



- Given a pair of plates:
  - For each activator, rank correlations (of markers), compute spearman correlation of ranking
  - Distance = 1-min correlation over activators

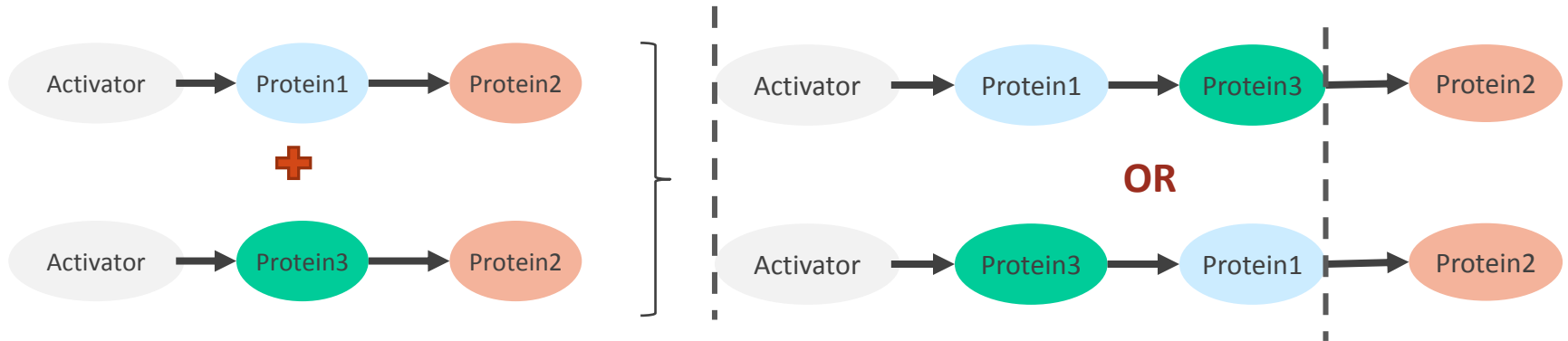Pipeline for making causal predictions

# Causal Postulates

| | | | | | |
|---|---|---|---|---|---|
| $PVO_4$ → pPlcg2 $\xrightarrow{0.5482}$ pSTAT3 | 0.875 | cd14-hladr- | | | |
| $PVO_4$ → pPlcg2 $\xrightarrow{0.5512}$ pZap70 | 0.8125 | cd14-hladrmid | | | |
| $PVO_4$ → pSlp76 $\xrightarrow{0.7152}$ pSHP2 | 0.8125 | cd14-hladrmid | | | |
| $PVO_4$ → pSHP2 $\xrightarrow{0.6708}$ pSTAT3 | 0.7857 | dendritic | | | |
| $PVO_4$ → pPlcg2 $\xrightarrow{0.8526}$ pP38 | 0.75 | cd14+hladr- | | | |
| $PVO_4$ → pPlcg2 $\xrightarrow{0.6166}$ pZap70 | 0.75 | cd14-hladr- | | | |
| $PVO_4$ → pSlp76 $\xrightarrow{0.5688}$ pZap70 | 0.75 | cd14-hladr- | | | |
| $PVO_4$ → pSHP2 $\xrightarrow{0.5688}$ pZap70 | 0.7143 | cd14-hladrmid | | | |
| $PVO_4$ → pSTAT3 $\xrightarrow{0.4557}$ pBtk | 0.7059 | cd14-hladr- | | | |
| BCR → pS6 $\xrightarrow{0.4557}$ pErk | 0.7037 | igm- | | | |

288 predictions in 14 sub-populations

- A list of **predicted causal** pairs, each "tagged" for a **specific population** and **activator,** ranked according to a score quantifying the frequency of appearance.
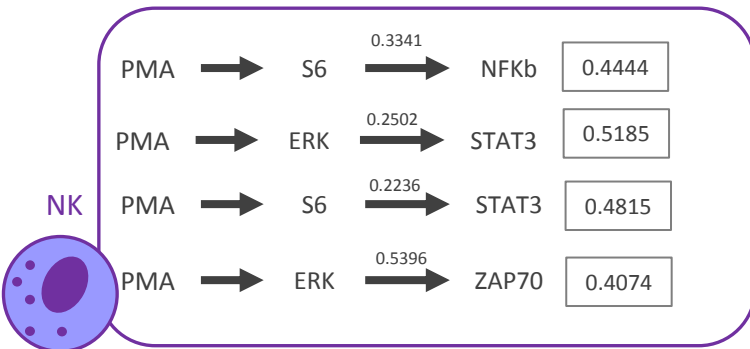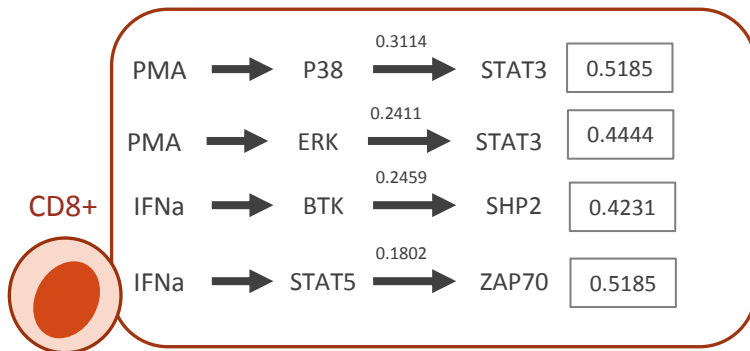
# Internal Validation



Check whether predicted
triplet has also been reported

- 42% of the predicted triplets are also reported
- Despite strict thresholds and multiple testing

- Theory+algorithms: [Tillman et. al. 2008, Triantafillou et. al 2010, Tsamardinos et. al 2012]
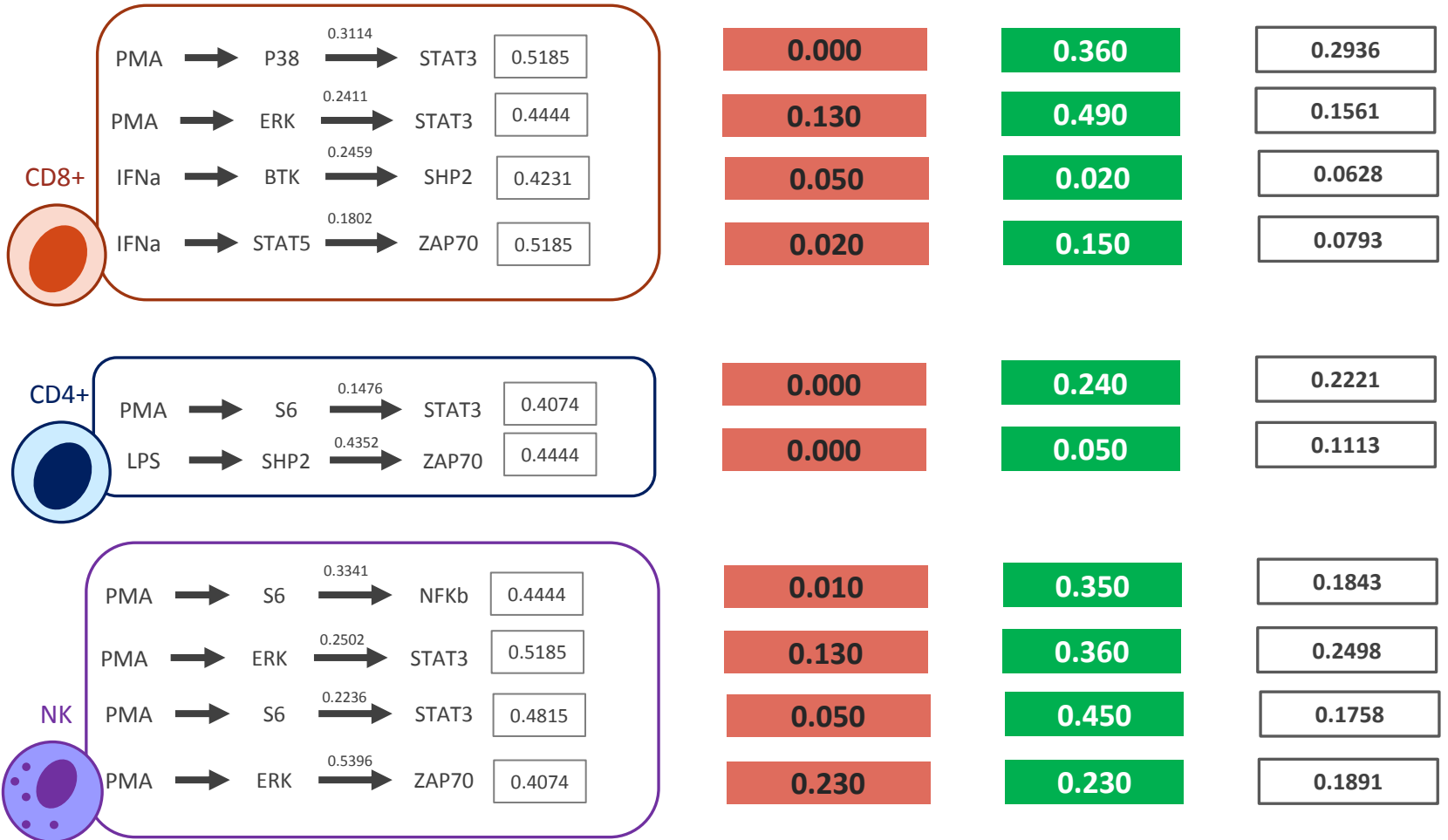
# Validation on Bendall Data



**CD8+**

PMA → P38 → $\overset{0.3114}{}$ STAT3   0.5185

PMA → ERK → $\overset{0.2411}{}$ STAT3   0.4444

IFNa → BTK → $\overset{0.2459}{}$ SHP2   0.4231

IFNa → STAT5 → $\overset{0.1802}{}$ ZAP70   0.5185

**CD4+**

PMA → S6 → $\overset{0.1476}{}$ STAT3   0.4074

LPS → SHP2 → $\overset{0.4352}{}$ ZAP70   0.4444

**NK**

PMA → S6 → $\overset{0.3341}{}$ NFKb   0.4444

PMA → ERK → $\overset{0.2502}{}$ STAT3   0.5185

PMA → S6 → $\overset{0.2236}{}$ STAT3   0.4815

PMA → ERK → $\overset{0.5396}{}$ ZAP70   0.4074

**Bendall Data**

- Run FCI with $a = 0.05$
- Bootstrap for robustness
- Report
  - **Conflicting** structures: Structures where $P_2 \rightarrow P_1$
  - **Confirming** Structures: Structures where $P_1 \rightarrow P_2$

⚠ Measurements in Bendall data are taken 15 minutes after activation

# Validation on Bendall Data



|  | | Conflicting | Confirming | Correlation |
|---|---|---|---|---|
| **CD8+** | PMA → P38 →(0.3114) STAT3 — 0.5185 | 0.000 | 0.360 | 0.2936 |
|  | PMA → ERK →(0.2411) STAT3 — 0.4444 | 0.130 | 0.490 | 0.1561 |
|  | IFNa → BTK →(0.2459) SHP2 — 0.4231 | 0.050 | 0.020 | 0.0628 |
|  | IFNa → STAT5 →(0.1802) ZAP70 — 0.5185 | 0.020 | 0.150 | 0.0793 |
| **CD4+** | PMA → S6 →(0.1476) STAT3 — 0.4074 | 0.000 | 0.240 | 0.2221 |
|  | LPS → SHP2 →(0.4352) ZAP70 — 0.4444 | 0.000 | 0.050 | 0.1113 |
| **NK** | PMA → S6 →(0.3341) NFKb — 0.4444 | 0.010 | 0.350 | 0.1843 |
|  | PMA → ERK →(0.2502) STAT3 — 0.5185 | 0.130 | 0.360 | 0.2498 |
|  | PMA → S6 →(0.2236) STAT3 — 0.4815 | 0.050 | 0.450 | 0.1758 |
|  | PMA → ERK →(0.5396) ZAP70 — 0.4074 | 0.230 | 0.230 | 0.1891 |

# Results

- Hundreds of predictions to-be-tested; Experiments under way!

- Internal validation using non-trivial inferences

- Promising validation on another collection of dataset (Bendall)

- Evidence of batch effects and/or biological reasons of variability

- Method based on the most basic causal discovery assumptions

# A Not So Basic Approach

# Co-analyzing data sets from different experimental conditions with overlapping variable sets



| p1 | p2 | ... | p30 |
|----|----|-----|-----|
|    |    |     |     |
|    |    |     |     |
|    |    |     |     |

Condition A

| p1 | p2 | ... | p30 |
|----|----|-----|-----|
|    |    |     |     |
|    |    |     |     |
|    |    |     |     |

Condition B

- Different experimental conditions
- Different variable sets

| p1 | p2 | ... | p30 |
|----|----|-----|-----|
|    |    |     |     |
|    |    |     |     |
|    |    |     |     |

Condition C

| p29 | p30 | ... | p40 |
|-----|-----|-----|-----|
|     |     |     |     |
|     |     |     |     |
|     |     |     |     |

Condition D

- Data can not be pulled together because they come from different distributions

- Principles of causality links them to the underlying causal graph

# Co-analyzing data sets from different experimental conditions with overlapping variable sets
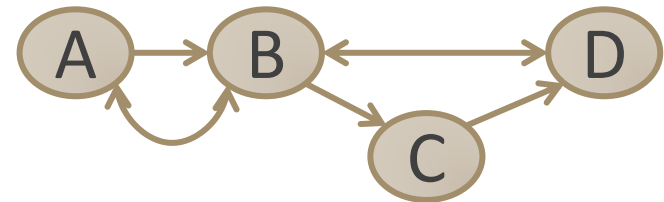


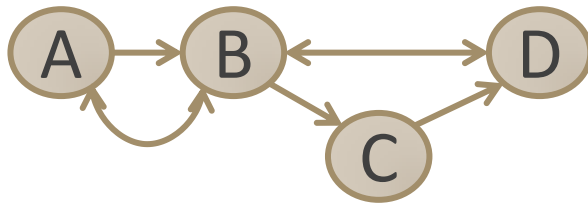Identify **a single causal** graph that simultaneously fits all data

# What type of causal graph?

- Semi-Markov causal models.

- $X \to Y$: $X$ causes $Y$ directly in the context of observed variables.

- $X \leftrightarrow Y$: $X$ and $Y$ share a latent common cause.

- Under faithfulness, $m$-separation entails all and only conditional independencies that stem from Causal Markov Condition.

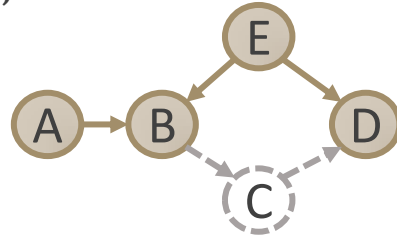- No learning algorithm.

# Manipulations in SMCMs



Graph (SMCM) $S$

- Values of $B$ are set solely by the manipulation procedure

- Graph surgery: Remove all edges into the manipulated node.
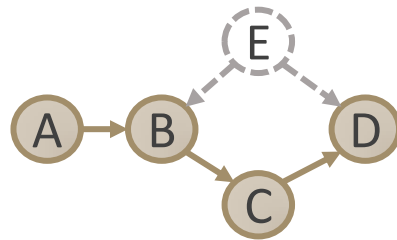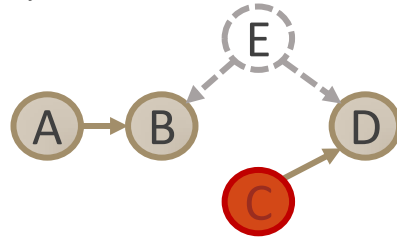
# Reverse Engineering

$S$, $C$ is latent



$Dep(A, D|\emptyset)_{D_1}$
$Dep(A, D|B)_{D_1}$
$Dep(A, D|E)_{D_1}$
$Ind(A, D|B, E)_{D_1}$
$Dep(A, B|\emptyset)_{D_1}$

...

$S$, $E$ is latent



$Dep(A, D|\emptyset)_{D_2}$
$Dep(A, D|B)_{D_2}$
$Dep(A, D|E)_{D_2}$
$Dep(A, D|C)_{D_2}$
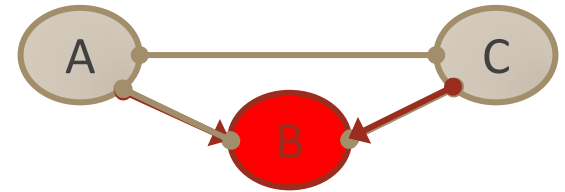$Dep(A, D|B, C)_{D_2}$

...

$S^C$, $E$ is latent



$Ind(A, D|\emptyset)_{D_3}$
$Dep(A, B|\emptyset)_{D_3}$
$Dep(A, B|C)_{D_3}$
$Dep(A, B|D)_{D_3}$
$Dep(A, D|Cd\ D)_{D_3}$
$Ind(B, C|\emptyset)_{D_3}$

...

$S$



Unknown True SMCM $S$

$S$ under manipulation and marginalization

Observed (in) dependencies

# Independencies as constraints

- Suppose you don't know anything about the structure $S$ of the three variables.

- You find out that in $S^B : Ind(A, C | \emptyset)$

- In path terms: $\nexists$ path in $S^B$ that is m-connecting $A$ and $C$ given $\emptyset$

- In SAT terms:

$$\neg edge(A, C) \wedge$$
$$[\neg edge(A, B) \vee arrow(A, B) \vee edge(B, \mathrm{C}) \vee arrow(C, B)]$$



A-C does not exist

AND

(A-B does not exist

OR

A-B is into B

OR

B-C does not exist

OR

B-C is into B)

# Statistical errors

- Constraints correspond to *
  1. Dependencies $Dep(A, B|\boldsymbol{Z})_{D_i}$
  2. Independencies $Ind(C, D|\boldsymbol{W})_{D_i}$

  ◦ **e.g.,** $Ind(A, B|\emptyset)_{D_1} \leftrightarrow \neg edge(A, C) \wedge [\neg edge(A, B) \vee arrow(A, B) \vee edge(B, C) \vee arrow(C, B)]$

- Compare a dependence to an independence
  ◦ How?
  ◦ Low p-value suggests dependence
  ◦ High p-value suggests independence
      (in the respective data set)



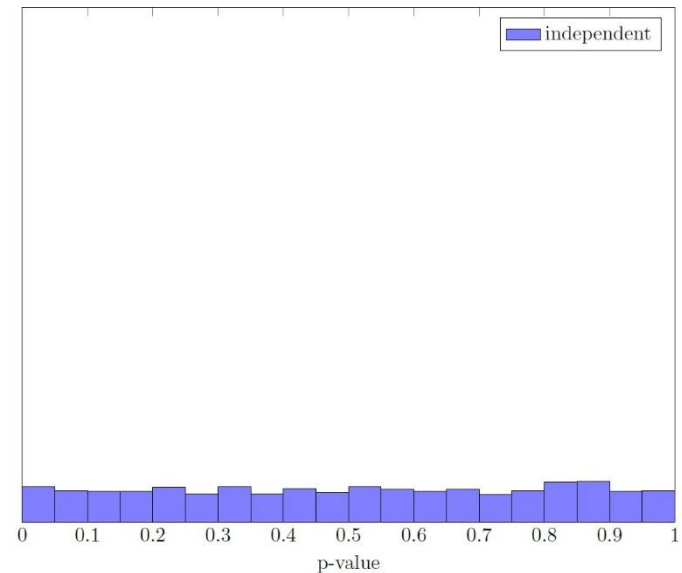Sort constraints!

What happens with statistical errors?

Conflicts make SAT instance unsatisfiable!

*well, not really

# Comparing p-values

- $H_0: p \sim Beta(1,1)$

- $H_1: p \sim Beta(\xi, 1), \ \xi \in (0,1)$

- $f(p|\pi_o, \xi) = \pi_0 + (1-\pi_0)\xi p^{\xi-1}$ , $\pi_0$: The proportion of p-values coming from $H_0$

- If you know $\widehat{\pi_0}, \ \hat{\xi}$ you can find the MAP ratio

- $E_0(p) = \frac{P(H_0|p)P(H_0)}{P(H_1|p)P(H_1)} = \frac{\widehat{\pi_0}}{(1-\widehat{\pi_0})\hat{\xi}p^{(1-\widehat{\xi})}}$, E$_1$ = 1/E$_0$

  ◦ If $E(p) > E(p)^{-1}$, independence is more likely
    than dependence

- **Sort p-values by max(E$_0$, E$_1$)**

- Use (Storey and Tibshirani, 2003) to identify $\widehat{\pi_o}$

- Minimize negative log likelihood of
  $f(p|\widehat{\pi_0}, \xi) = \widehat{\pi_0} + (1-\widehat{\pi_0})\xi p^{\xi-1}$ to identify $\hat{\xi}$ .

- Rank constraints according to MAP ratio and satisfy them if
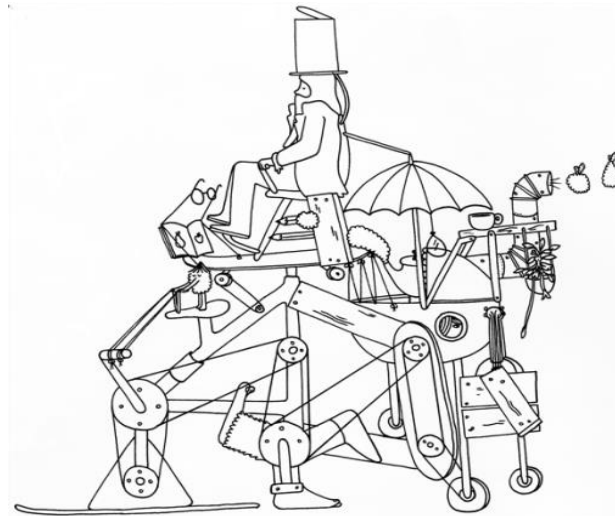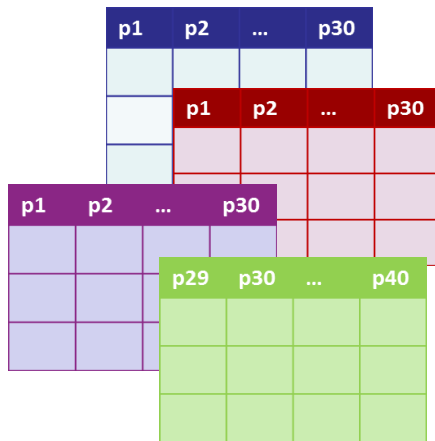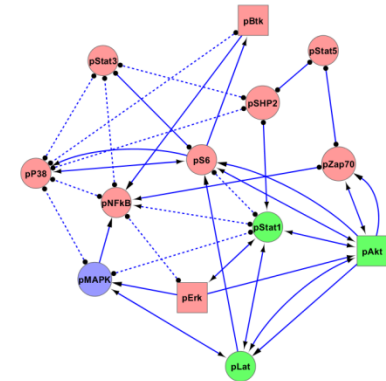  possible in the given order.

# "COmbINE" Algorithm

Data sets $D_i$ measuring overlapping variables under different experimental conditions

COmbINE
Algorithm that transforms independence constrains to SAT instance

Summary of semi Markov Causal models that best fits all data sets simultaneously



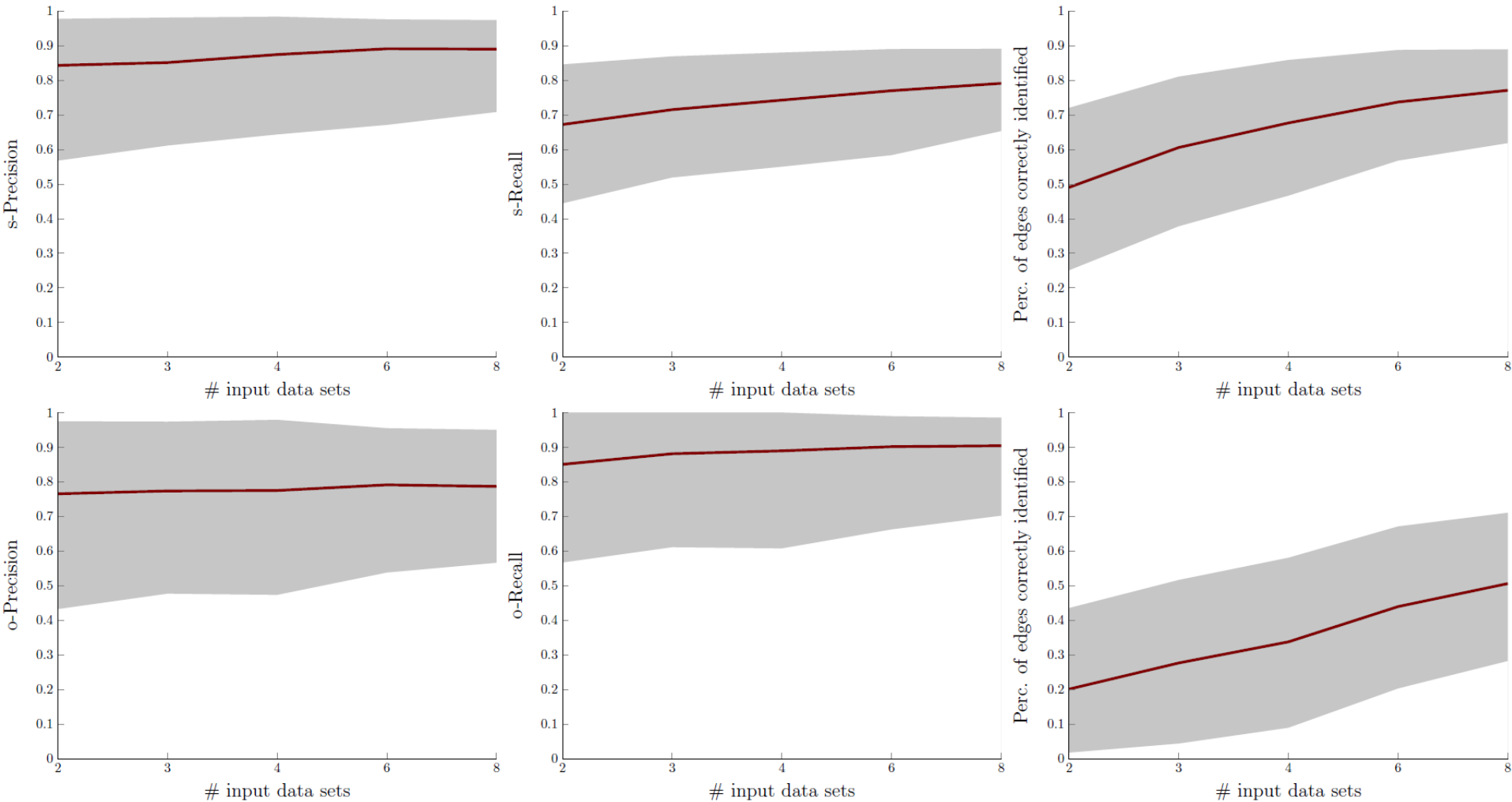Eric Ellis

# Similar Algorithms

- SBCSD: [Hyttinen et al., UAI, 2013]
  - Inherently less compact representation of path constraints.
  - Does not handle conflicts; non applicable to real data.
  - In addition, it admits cycles.
  - Scales up to 14 variables

- Lininf  [Hyttinen et al., UAI 2012, JMLR 2012]
  - Linear relations only.
  - Scales up poorly (6 variables in total with overlapping variables, 10 without).
  - In addition, it admits cycles.

|      | COmbINE              | SBCSD                  |
|------|----------------------|------------------------|
| ASIA | $7.1768 \pm 5.2424$  | $51.6617 \pm 27.5997$  |
| CAR  | $3.6994 \pm 2.2489$  | $211.5117 \pm 78.2334$ |

Execution Time in Seconds

# Performance on Simulated Data

# Application on Mass Cytometry data



cd4+ T-cells

cd8+ T-cells

Response to PMA

| Data set | Source | $\mathbf{L_i}$ | $\mathbf{I_i}$ | Donor |
|:---:|:---:|:---:|:---:|:---:|
| $\mathbf{D_1}$ | Bodenmiller et al. (2012) | pMAPK | pAkt | 1 |
| $\mathbf{D_2}$ | Bodenmiller et al. (2012) | pMAPK | pBtk | 1 |
| $\mathbf{D_3}$ | Bodenmiller et al. (2012) | pMAPK | pErk | 1 |
| $\mathbf{D_4}$ | Bendall et al. (2011) | pAkt, pLat, pStat1 | pErk | 2 |
| $\mathbf{D_5}$ | Bendall et al. (2011) | pAkt, pLat, pStat1 | pErk | 3 |

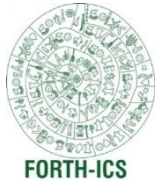# Summary and Conclusions

- Mass Cytometry data a good domain for causal discovery

- Hundreds of robust causal postulates

- Approach:
  ◦ Conservative: local discovery, performing all tests, independent analysis of populations
  ◦ Opportunistic: using 2 thresholds for (in)dependency

- New algorithm that can handle
  ◦ different experimental conditions
  ◦ overlapping variable subsets
  ◦ deal with statistical errors

- Numerous directions open for future work on this collection of data
  ◦ Experiments under way!

# Acknowledgements and Credit

**Ioannis Tsamardinos**
Associate Prof
Lab Head

**Jesper Tegnér**
Prof
Unit Head

**Sofia Triantafillou**
Ph.D. Candidate

**Angelika Schmidt**
Post-Doc

**Vincenzo Lagani**
Research Fellow

**David Gomez-Cabrero,**
Project Leader