# Reactmine: an algorithm for inferring biochemical reactions from time series data

Julien Martinelli, Jeremy Grignard, Sylvain Soliman, Annabelle Ballesta, François Fages
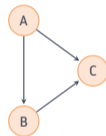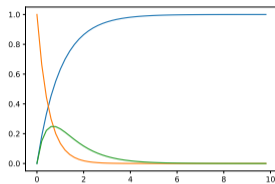
Tuesday, January 3rd 2023

# Network Inference from time-series data

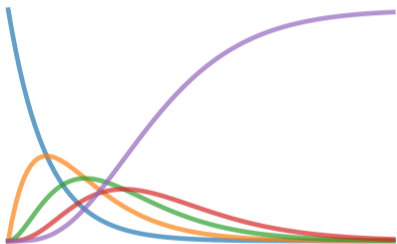Input: time series describing evolution of molecular species

Output: interaction graph

- (un)oriented graph → Gene Regulatory Network inference
  - ▸ Gaussian Processes   *(Aalto et al. 2019)*
  - ▸ Information theory   *(Chan et al. 2017)*
  - ▸ Correlation networks   *(Krumsiek et al. 2011)*
- oriented / weighted graph → **Chemical Reaction Network Inference**
  - ▸ Evolutionary algorithms   *(Choi et al. 2018)*
  - ▸ Sparse regression   *(Brunton et al. 2016)*

# Chemical Reaction Network Inference

Input: single time series data $Y = \left(y_{l,i}\right)_{\substack{1 \leqslant l \leqslant n \\ 1 \leqslant i \leqslant m}}$

Output:
Chemical Reaction Network



| Hidden CRN | Learned CRN |
|:---:|:---:|
| $A \xrightarrow{1} B$ | $A \xrightarrow{0.999} B$ |
| $B \xrightarrow{1} C$ | $B \xrightarrow{1.001} C$ |
| $C \xrightarrow{1} D$ | $C \xrightarrow{1.002} D$ |
| $D \xrightarrow{1} E$ | $D \xrightarrow{0.999} E$ |

$A$ → 0.999 → $B$ → 1.001 → $C$ → 1.002 → $D$ → 0.999 → $E$

For this presentation: $A \xrightarrow{k} B \iff \begin{cases} \dot{A} = -kA \\ \dot{B} = kA \end{cases}$

# Framework

Reaction: $(R, P, f)$ with $R$ (resp. $P$) set of reactants (resp. products) and $f$ rate function.

Chemical Reaction Network (CRN): Finite set of reactions

# Framework

Reaction: $(R, P, f)$ with $R$ (resp. $P$) set of reactants (resp. products) and $f$ rate function.

Chemical Reaction Network (CRN): Finite set of reactions

- 0/1 Stoichiometry

- Elementary reactions: at most two reactants

- At most 1 catalyst (e.g. $D$ in $A + D \xrightarrow{k} B + D$)

# Framework

Reaction: $(R, P, f)$ with $R$ (resp. $P$) set of reactants (resp. products) and $f$ rate function.

Chemical Reaction Network (CRN): Finite set of reactions

- 0/1 Stoichiometry

- Elementary reactions: at most two reactants

- At most 1 catalyst (e.g. $D$ in $A + D \xrightarrow{k} B + D$)

Learning protocol:

▶ Learn a CRN involving only observed species

▶ Based on a single trace (no combinatorics of initial states and *knockouts*)

# Backbone of most methods: Sparse Identification of Nonlinear Dynamics

$$\Xi = \underset{\Xi \in \mathbb{R}^{p \times m}}{\operatorname{argmin}} \|\dot{Y} - \Theta(Y)\Xi\|_F^2 + \lambda\|\Xi\|_1$$

$\Theta(Y) \in \mathbb{R}^{n \times p}$: library of $p$ functions, e.g.
$$\begin{bmatrix} | & | & & | & | & & | \\ 1 & Y_{\bullet,1} & \dots & Y_{\bullet,m} & Y_{\bullet,1}Y_{\bullet,2} & \dots & Y_{\bullet,m-1}Y_{\bullet,m} \\ | & | & & | & | & & | \end{bmatrix}$$
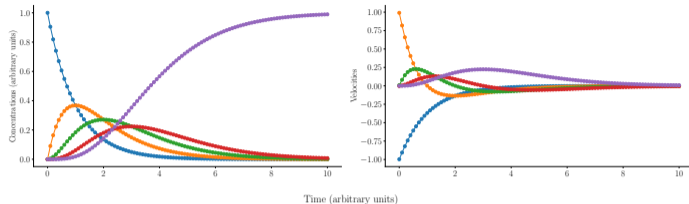
$Y_{\bullet,i}$: time concentration vector for species $i$

$\Xi \in \mathbb{R}^{p \times m}$: weight matrix

$\lambda \in \mathbb{R}^+$: hyperparameter controlling level of sparsity

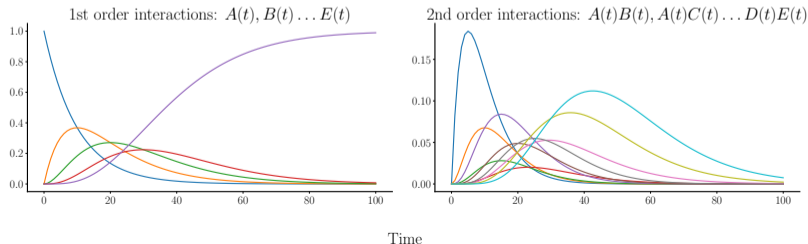# Example - Chain Chemical Reaction Network

$$\begin{cases} \dot{A} = -A \\ \dot{B} = A - B \\ \dot{C} = B - C \\ \dot{D} = C - D \\ \dot{E} = D \end{cases}$$



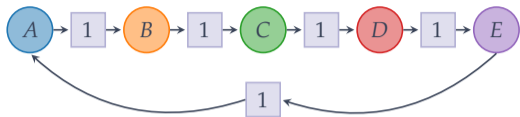Left/right plots: simulated concentrations, derivatives
SINDy aims to predict the derivatives with the following library functions

Library members - Chain chemical reaction network



1st order interactions: $A(t), B(t) \ldots E(t)$

2nd order interactions: $A(t)B(t), A(t)C(t) \ldots D(t)E(t)$

Time

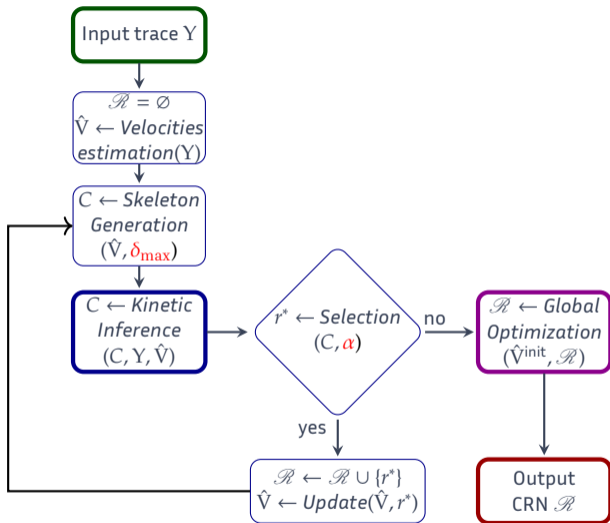# SINDy fails at CRN inference



$$\begin{cases} \dot{A} = E - A \\ \dot{B} = A - B \\ \dot{C} = B - C \\ \dot{D} = C - D \\ \dot{E} = D - E \end{cases}$$
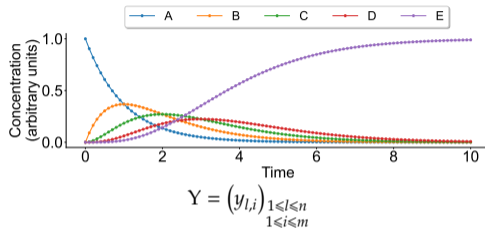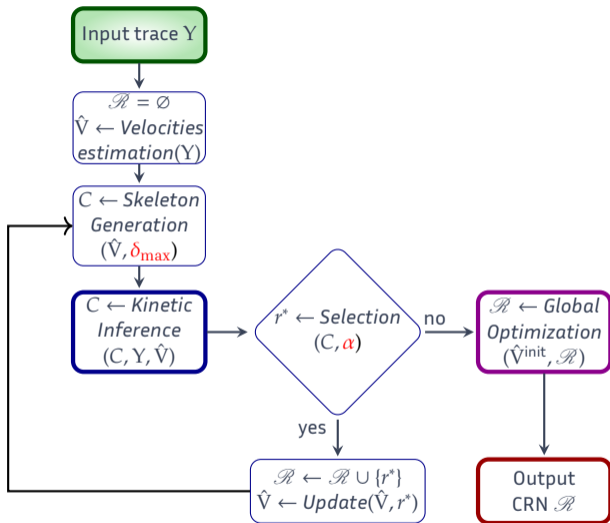
$$\begin{cases} \dot{A} = -1.00A + 1.03E - 0.006D - 0.07AE - 0.06DE \\ \dot{B} = 1.00A - 1.00B + 0.004C + 0.001AB - 0.211AC - 0.092BC \\ \dot{C} = 1.14B - 1.18C - 0.002D - 0.17AB + 0.39CD \\ \dot{D} = 0.35B - 0.35E \\ \dot{E} = 0.39C + 0.457E - 4.21AE \end{cases}$$

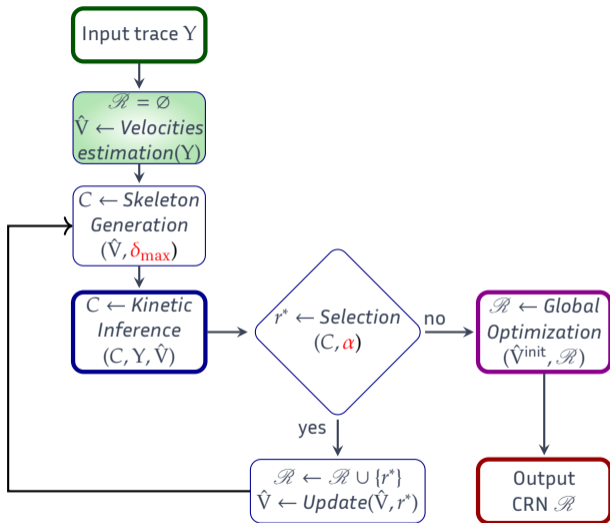Best ODE system found across all sparsity threshold $\lambda$

# Core Reactmine sequential algorithm
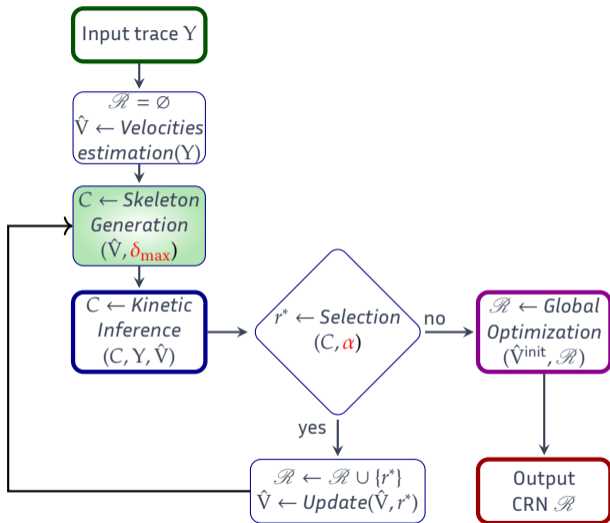
# Core Reactmine sequential algorithm



$$Y = \left( y_{l,i} \right)_{\substack{1 \leqslant l \leqslant n \\ 1 \leqslant i \leqslant m}}$$

# Core Reactmine sequential algorithm



$$\text{Velocities } \hat{V} = \left(\hat{v}_{l,i}\right)_{\substack{1 \leqslant l \leqslant n \\ 1 \leqslant i \leqslant m}}$$

$$\hat{v}_{l,i} = \frac{y_{l+1,i} - y_{l,i}}{t_{l+1} - t_l}$$
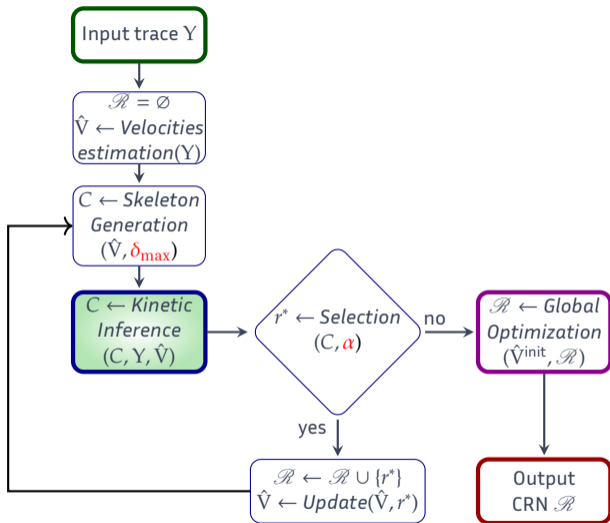
# Core Reactmine sequential algorithm



Each reaction skeleton $r = (R, P)$
is inferred based on time points $t_l$
where it is **preponderant: support set** $\mathcal{T}(r)$

Reactants and products belonging to a skeleton
have similar absolute velocities up to $\delta_{max}$

# Core Reactmine sequential algorithm



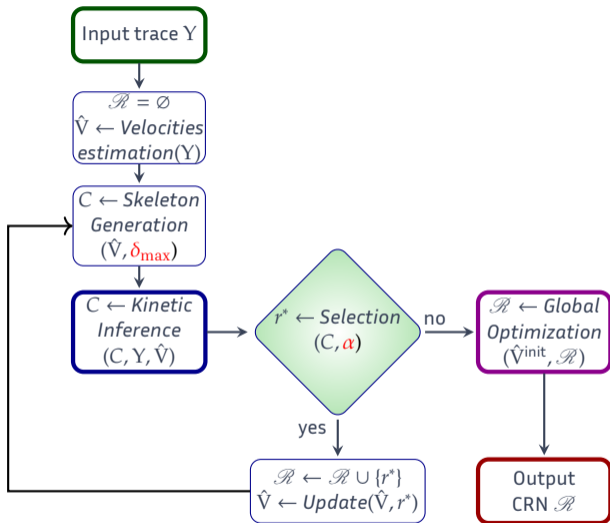For each reaction skeleton $r = (R, P)$ associate kinetic rate

$$\forall j \in R \cup P, \ \forall l \in \{1, \ldots, n\}, \ |v_{l,j}| = k \prod_{u \in R} y_{l,u}$$

**Estimate k reliably on the support set $\mathscr{T}(r)$**

$$\hat{k} = \frac{1}{\#\mathscr{T}(r)} \sum_{l \in \mathscr{T}(r)} \frac{|\hat{v}_{l,j}|}{\prod_{u \in R} y_{l,u}}$$

Coefficient of variation (CV) $\rho = \frac{\sigma}{|\hat{k}|}$

8

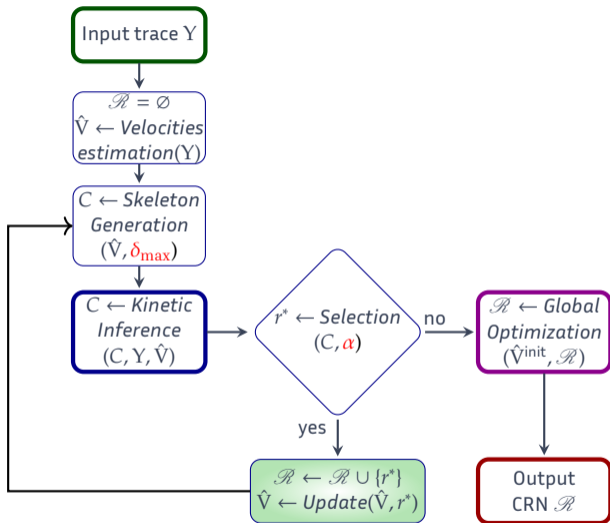# Core Reactmine sequential algorithm



**Select** reaction minimizing CV
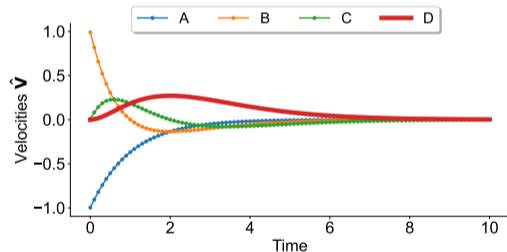$$r^* = \operatorname*{argmin}_{r} \rho_r$$

**Accept** $r^*$ if $\rho_{r^*} < \alpha$

# Core Reactmine sequential algorithm



Remove the effect of accepted reaction on the velocities

$$\hat{V} \leftarrow \hat{V} - \begin{pmatrix} f(Y_{1,\bullet}) \\ \vdots \\ f(Y_{n,\bullet}) \end{pmatrix} s^T$$

effect of the reaction          stoichiometry vector

$Y_{l,\bullet}$: species concentration vector at time $t_l$

# Core Reactmine sequential algorithm



Joint optimization
of kinetic parameters
**over whole trace**

$$k = \underset{k \in \mathbb{R}_+^p}{\mathrm{argmin}} \| \hat{V}^{init} - F(Y,k)S \|_F^2$$

$= \Delta$ whole trace CRN transition discrepancy

# Core Reactmine sequential algorithm



Input trace $Y$

$\mathcal{R} = \varnothing$
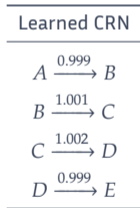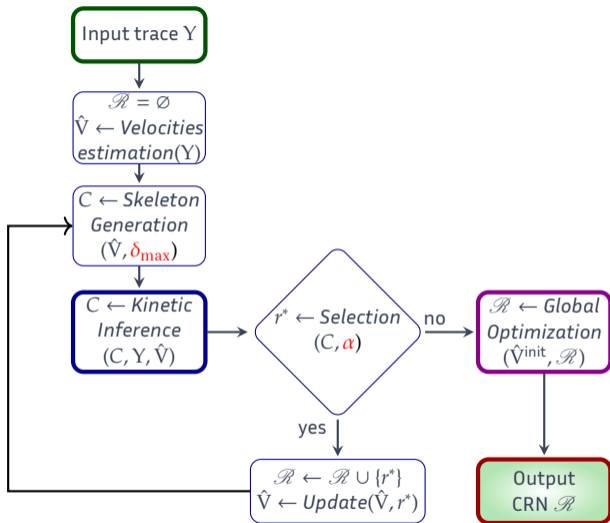$\hat{V} \leftarrow$ *Velocities estimation*$(Y)$

$C \leftarrow$ *Skeleton Generation* $(\hat{V}, \delta_{max})$

$C \leftarrow$ *Kinetic Inference* $(C, Y, \hat{V})$

$r^* \leftarrow$ *Selection* $(C, \alpha)$

no

yes

$\mathcal{R} \leftarrow$ *Global Optimization* $(\hat{V}^{init}, \mathcal{R})$

Output CRN $\mathcal{R}$

$\mathcal{R} \leftarrow \mathcal{R} \cup \{r^*\}$
$\hat{V} \leftarrow$ *Update*$(\hat{V}, r^*)$

Learned CRN

$A \xrightarrow{0.999} B$

$B \xrightarrow{1.001} C$

$C \xrightarrow{1.002} D$

$D \xrightarrow{0.999} E$

$A \xrightarrow{0.999} B \xrightarrow{1.001} C \xrightarrow{1.002} D \xrightarrow{0.999} E$

# Example - iterative inference of reactions for Chain CRN



Support of reaction MA(0.99869) for D → E

**Example - iterative inference of reactions for Chain CRN**
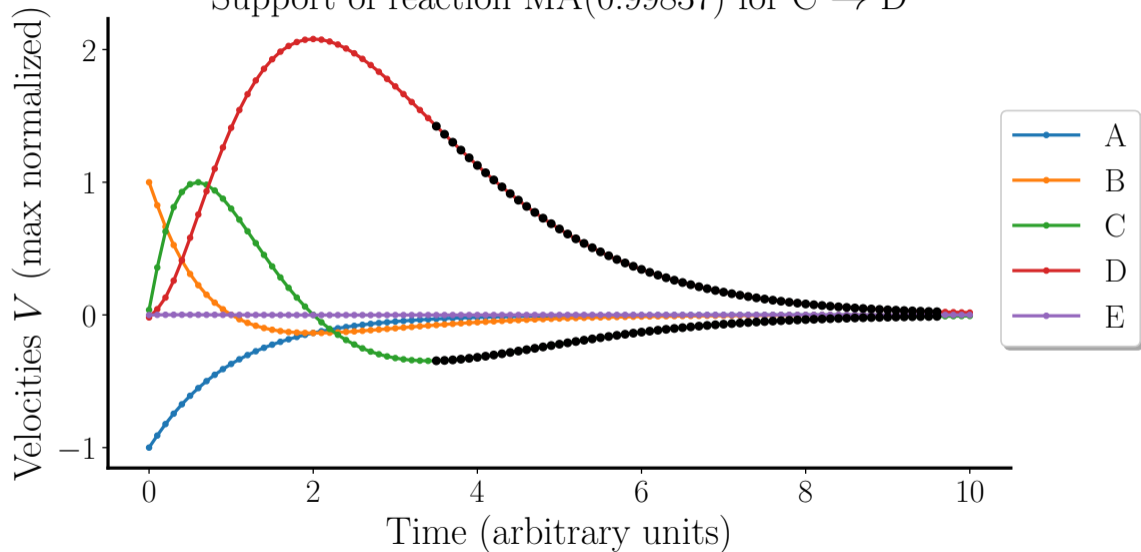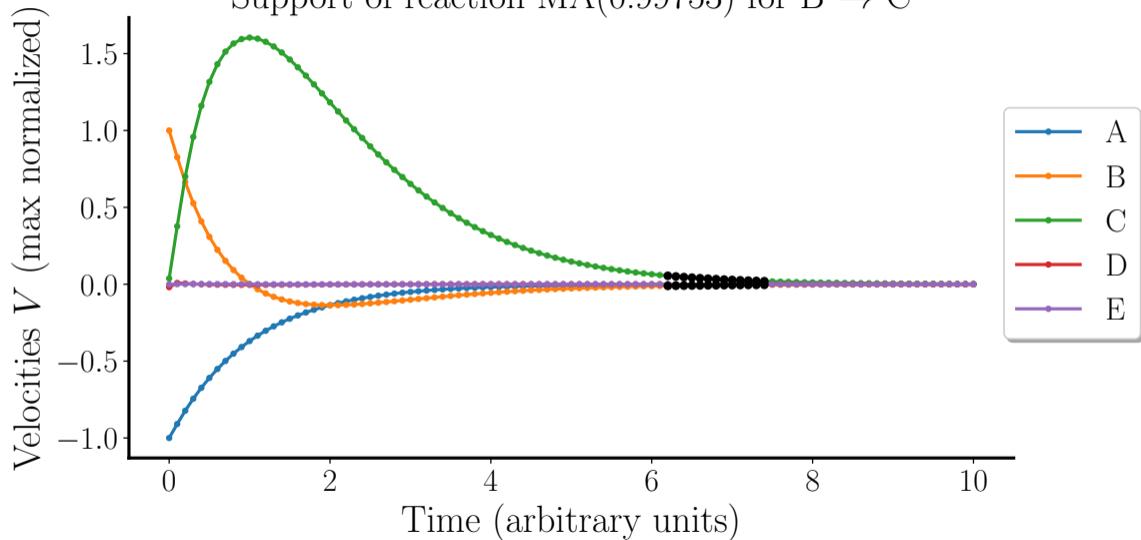
Support of reaction MA(0.99837) for C → D

# Example - iterative inference of reactions for Chain CRN



Support of reaction MA(0.99753) for B → C

# Example - iterative inference of reactions for Chain CRN

Support of reaction MA(1.00069) for A → B

# Reactmine search algorithm
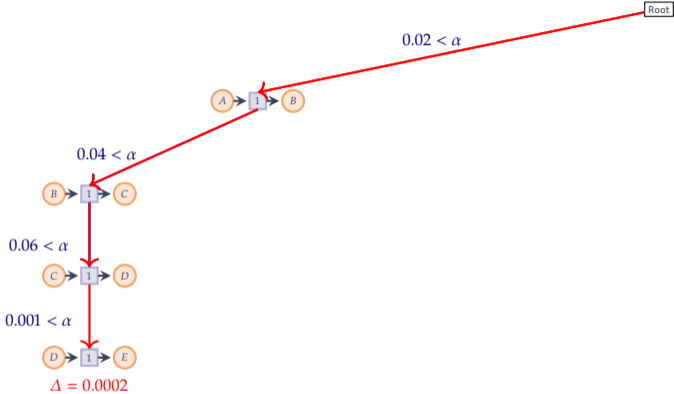
# Reactmine search algorithm



number of candidates per node: $\beta = 3$

node accepted if quality measure $\rho < \alpha$

$0.02 < \alpha$

$0.02 < \alpha$

$0.1 < \alpha$

$0.04 < \alpha$

$0.06 < \alpha$

$0.001 < \alpha$

$\Delta = 0.0002$

# Reactmine search algorithm



number of candidates per node: $\beta = 3$

node accepted if quality measure $\rho < \alpha$

$0.02 < \alpha$

$0.02 < \alpha$

$0.1 < \alpha$

$0.04 < \alpha$

$0.5 > \alpha$

$0.06 < \alpha$

$0.001 < \alpha$

$\Delta = 0.0002$

# Reactmine search algorithm



number of candidates per node: $\beta = 3$

node accepted if quality measure $\rho < \alpha$

**Termination:** $\rho > \alpha$ for all reactions
or reached $\gamma$ maximal depth allowed

$0.02 < \alpha$

$0.02 < \alpha$

$0.1 < \alpha$
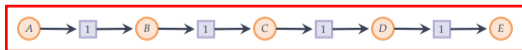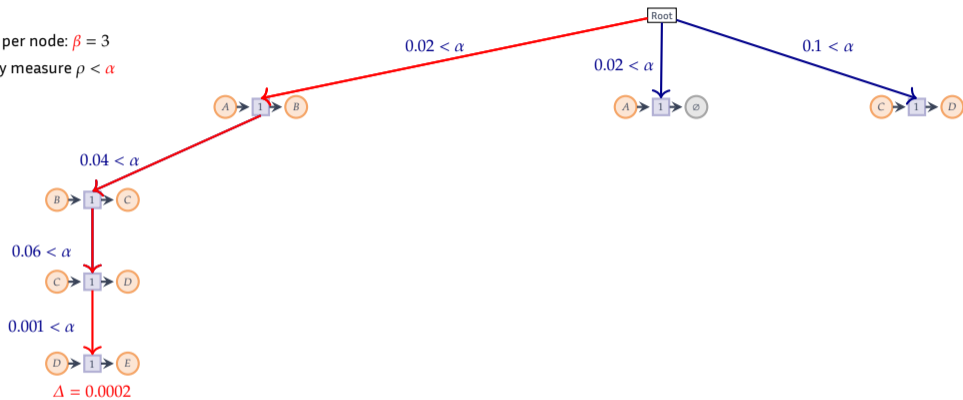
$0.04 < \alpha$

$0.5 > \alpha$

$0.06 < \alpha$

$0.001 < \alpha$

$\Delta = 0.0002$

# Reactmine search algorithm

number of candidates per node: $\beta = 3$

node accepted if quality measure $\rho < \alpha$

**Termination:** $\rho > \alpha$ for all reactions
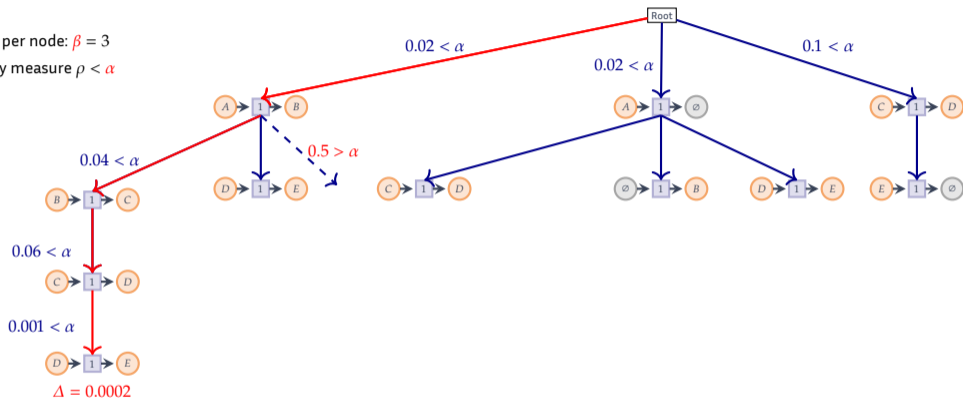or reached $\gamma$ maximal depth allowed

# Reactmine search algorithm

number of candidates per node: $\beta = 3$

node accepted if quality measure $\rho < \alpha$

**Termination:** $\rho > \alpha$ for all reactions
or reached $\gamma$ maximal depth allowed
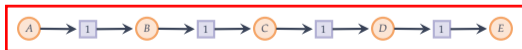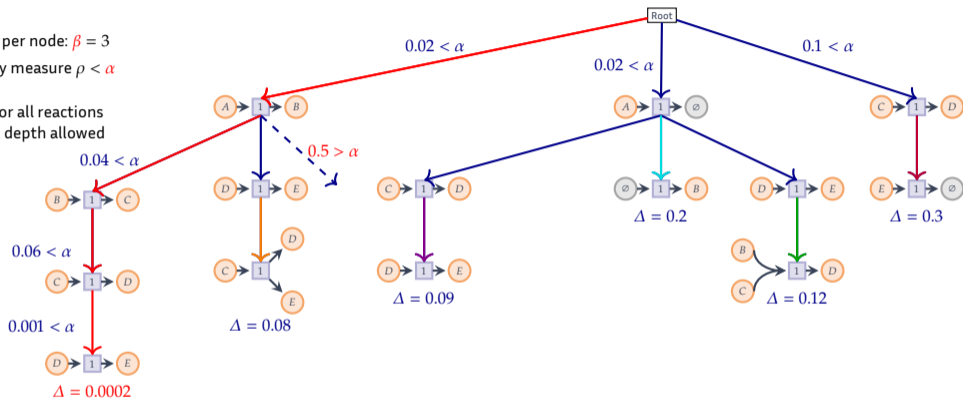


10

# Reactmine search algorithm

number of candidates per node: $\beta = 3$

node accepted if quality measure $\rho < \alpha$

**Termination:** $\rho > \alpha$ for all reactions
or reached $\gamma$ maximal depth allowed



4 Hyperparameters:

- $\delta_{max}$ Species variations similarity threshold
- $\alpha$ CV threshold
- $\gamma$ CRN size limit
- $\beta$ Number of reaction candidates per node

# Reactmine search algorithm

number of candidates per node: $\beta = 3$
node accepted if quality measure $\rho < \alpha$

**Termination:** $\rho > \alpha$ for all reactions
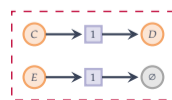or reached $\gamma$ maximal depth allowed
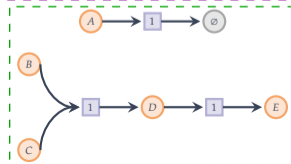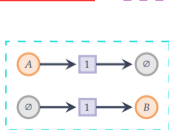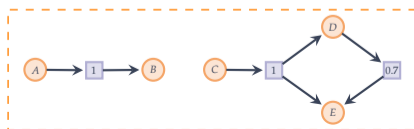


4 Hyperparameters:

- $\delta_{max}$ Species variations similarity threshold
- $\alpha$ CV threshold
- $\gamma$ CRN size limit
- $\beta$ Number of reaction candidates per node

**Hyperparameter selection by minimization of $\Delta$**

# Evaluation on Loop CRN



| Hidden CRN | Learned CRN |
|:---:|:---:|
| $A \xrightarrow{1} B$ | $A \xrightarrow{1} B$ |
| $B \xrightarrow{1} C$ | $B \xrightarrow{1} C$ |
| $C \xrightarrow{1} D$ | $C \xrightarrow{1} D$ |
| $D \xrightarrow{1} E$ | $D \xrightarrow{1} E$ |
| $E \xrightarrow{1} A$ | $E \xrightarrow{1} A$ |

Inference difficult → each species takes part in two reactions

$$\begin{cases} \dfrac{dA}{dt} = k_5 E - k_1 A \\[2mm] \dfrac{dB}{dt} = k_1 A - k_2 B \\[2mm] \dfrac{dC}{dt} = k_2 B - k_3 C \\[2mm] \dfrac{dD}{dt} = k_3 C - k_4 D \\[2mm] \dfrac{dE}{dt} = k_4 D - k_5 E \end{cases}$$

# Lokta-Volterra

| Ground-truth | Learned CRN | SINDy ODE |
|---|---|---|
| $A \xrightarrow{0.3} \varnothing$ | $A \xrightarrow{0.298} \varnothing$ | |
| $B \xrightarrow{1} 2B$ | $B \xrightarrow{0.994} 2B$ | $\begin{cases} \dot{A} = -0.299A + 0.010AB \\ \dot{B} = 0.995B - 0.010AB \end{cases}$ |
| $A + B \xrightarrow{0.01} 2A$ | $A + B \xrightarrow{0.010} 2A$ | |

# Simplified MAPK Cascade

$$KKK \xrightarrow{0.0045} KKKp$$

$$KKKp + KK \xrightarrow{1000} KKKpKK$$

$$KKKpKK \xrightarrow{150} KKKp + KK$$

$$KKKpKK \xrightarrow{150} KKKp + KKp$$

$$KKKp + KKp \xrightarrow{1000} KKKpKKp$$

$$KKKpKKp \xrightarrow{150} KKKp + KKp$$

$$KKKpKKp \xrightarrow{150} KKKp + KKpp$$

$$KKK \xrightarrow{0.0045} KKKp$$

$$KKp + KKKp \xrightarrow{499.97} KKpp + KKKp$$

$$KK + KKKp \xrightarrow{501.01} KKKp$$

$$KKKpKK \xrightarrow{150.04} KKKp + KKp$$

$$KKKp + KK \xrightarrow{501.19} KKKpKK$$

$$KKKpKK \xrightarrow{150.37} KKKp + KK$$

$$KKKp + KKp \xrightarrow{517.78} KKKpKKp$$

$$KKKpKKp \xrightarrow{155.34} KKKp + KKp$$

$$KKKp + KK \xrightarrow{500.19} KKKpKK + KK$$

# Simplified MAPK Cascade

$$KKK \xrightarrow{0.0045} KKKp$$

$$KKKp + KK \xrightarrow{1000} KKKpKK$$

$$KKKpKK \xrightarrow{150} KKKp + KK$$

$$KKKpKK \xrightarrow{150} KKKp + KKp$$

$$KKKp + KKp \xrightarrow{1000} KKKpKKp$$

$$KKKpKKp \xrightarrow{150} KKKp + KKp$$

$$KKKpKKp \xrightarrow{150} KKKp + KKpp$$

$$KKK \xrightarrow{0.0045} KKKp$$

$$KKp + KKKp \xrightarrow{499.97} KKpp + KKKp$$

$$KK + KKKp \xrightarrow{501.01} KKKp$$

$$KKKpKK \xrightarrow{150.04} KKKp + KKp$$

$$KKKp + KK \xrightarrow{501.19} KKKpKK$$

$$KKKpKK \xrightarrow{150.37} KKKp + KK$$

$$KKKp + KKp \xrightarrow{517.78} KKKpKKp$$

$$KKKpKKp \xrightarrow{155.34} KKKp + KKp$$

$$KKKp + KK \xrightarrow{500.19} KKKpKK + KK$$



- Adding up Reaction 3 and 9 is ODE-equivalent to $KKKp + KK \rightarrow KKKpKK$

- This reaction has already been inferred ($5^{th}$) $\rightarrow$ simplification

- The "simplified" inferred CRN has 7 reactions 6 of which are accurate.

## SINDy inferred ODE system for MAPK

$$\dot{KKpp} = 11764.89 - 9818.81 KKpp - 21809.21 KKKp - 64774.82 KKKpKKp - 9881.63 KKp$$
$$+109653.06 KKKpKK - 10102.57 KK - 23028.83 KKK + 23087.90 KKpp \times KKKp$$
$$+47383.24 KKpp \times KKKpKKp + 0.01 KKpp \times KKp - 94598.54 KKpp \times KKKpKK$$
$$+0.05 KKpp \times KK + 24104.25 KKpp \times KKK + 58119.14 KKKp \times KKp$$
$$+117674.08 KKKp \times KK + 68314.42 KKKpKKp \times KKp + 239788.36 KKKpKKp \times KK$$
$$-171491.97 KKp \times KKKpKK + 0.03 KKp \times KK + 45027.82 KKp \times KKK + 118690.34 KK \times KKK$$

$$\dot{KKKp} = 0.003 KKp - 0.001 KKpp \times KKp - 0.004 KKpp \times KK - 0.002 KKp \times KK$$

$$\dot{KKKpKKp} = -0.002 KKp + 0.001 KKpp \times KKp + 0.004 KKpp \times KK + 0.002 KKp \times KK$$

$$\dot{KKp} = -11345.814 + 9469.53 KKpp + 20253.98 KKKp + 62939.32 KKKpKKp$$
$$+9525.97 KKp - 105529.77 KKKpKK + 9401.39 KK + 22299.11 KKK$$
$$-21772.25 KKpp \times KKKp - 45730.13 KKpp \times KKKpKKp - 0.01 KKpp \times KKp$$
$$+91003.53 KKpp \times KKKpKK - 0.05 KKpp \times KK - 23476.51 KKpp \times KKK$$
$$-56249.34 KKKp \times KKp + 996.55 KKKp \times KK - 64537.43 KKKpKKp \times KKp$$
$$-113908.83 KKKpKKp \times KK + 163092.38 KKp \times KKKpKK - 0.03 KKp \times KK$$
$$-42276.50 KKp \times KKK + 113704.11 KKKpKK \times KK - 763.549 KK \times KKK$$
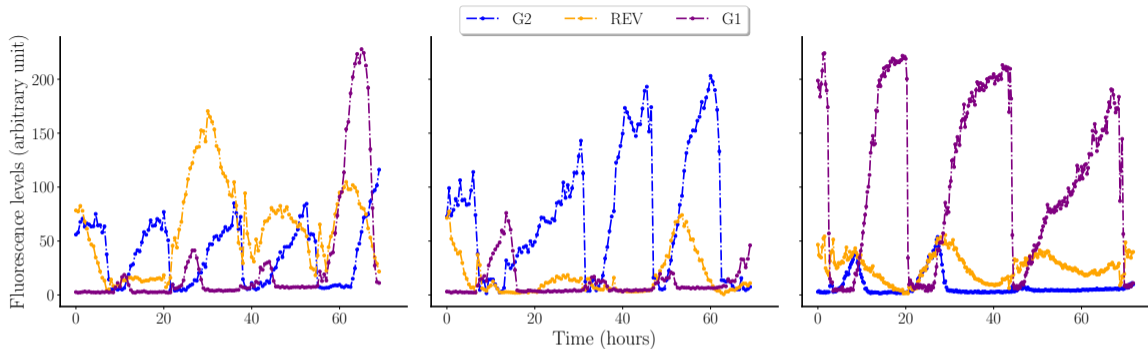
$$\dot{KKKpKK} = 0$$

$$\dot{KK} = -668.78 + 557.83 KKpp + 1724.07 KKKp + 2552.29 KKKpKKp$$
$$+564.79 KKp - 5063.51 KKKpKK + 551.20 KK + 784.94 KKK$$
$$-1609.99 KKpp \times KKKp + -1960.37 KKpp \times KKKpKKp - 0.001 KKpp \times KKp$$
$$+4399.19 KKpp \times KKKpKK - 0.01 KKpp \times KK - 827.38 KKpp \times KKK$$
$$-3441.84 KKKp \times KKp + 543.68 KKKp \times KK - 4279.75 KKKpKKp \times KKp$$
$$-8537.02 KKKpKKp \times KK + 10869.45 KKp \times KKKpKK - 0.003 KKp \times KK$$
$$-3146.138 KKp \times KKK + 6610.523 KKKpKK \times KK + 1384.463 KK \times KKK$$

$$\dot{KKK} = 0$$

# Application on real data: videomicroscopy

- NIH3T3 embryonic mouse fibroblasts left to proliferate in regular medium supplemented with 20% FBS concentration

- Time lapse videomicroscopy, one image taken every 15 minutes during 72 hours

- Cell tracking using three different fluorescent markers of the circadian clock and the cell cycle:
  - Reverb$\alpha$ circadian clock protein reporter
  - Fluorescence Ubiquitination Cell Cycle Indicators, Cdt1 and Geminin, two cell cycle proteins which accumulate during the G1 and S/G2/M phases, respectively.
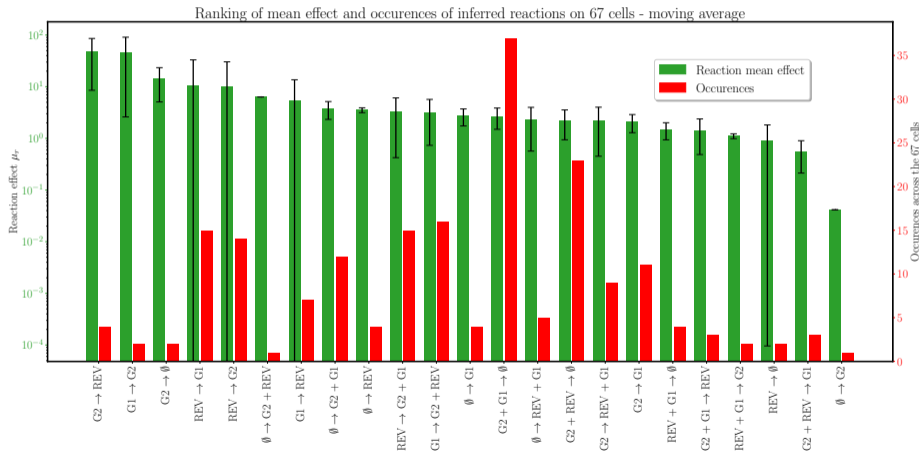
# Highly heterogeneous cell behavior



- 67 cells after curation
- Data smoothing using a moving average
- High heterogeneity → infer one CRN per cell
- We search for Michaelis-Menten reactions: $f(y) := v_{\max} \dfrac{y}{K_m + y}$

16

# Distribution of inferred reactions

1. For each cell, select the best CRN inferred
2. Compute $n_{occurrences}$ of reaction $r = (R, P, f)$ and mean effect $\mu_r = \frac{1}{nC} \sum_{c=1}^{C} \sum_{l=1}^{n} f(y_l^{(c)})$



Ranking of mean effect and occurences of inferred reactions on 67 cells - moving average

$G2 \rightarrow REV$ and $G1 \rightarrow G2$ recovered and present in literature

# Conclusion

- A method to **sequentially** infer biochemical reactions.

    ▸ **Parsimony** of the inferred network integrated by construction.

- Philosophy: **"mining" reactions** at specific time points where they are preponderant.

    ▸ More reliable estimation of reaction kinetics based on support

    ▸ **Explainability** of the method through the support set of inferred reaction

- Successfully tackled multi-scaled / cyclic CRNs

# Short and long term perspectives

Noise and real data:

- Proper evaluation/treatment against noisy data (e.g. bootstrap)
- Variables non observed at the same time points
- Non uniform grid of observation time points

# Short and long term perspectives

Noise and real data:

- Proper evaluation/treatment against noisy data (e.g. bootstrap)
- Variables non observed at the same time points
- Non uniform grid of observation time points

**Hidden species:**

- Assume two species $A, B$ for which $A(0)$ and $B(0)$ are available, but we only have the time series $X(t) = A(t) + B(t)$
  $\rightarrow$ Can we still infer a network involving $A$ and $B$?
- Infer completely unobserved hidden species?
- **→Evolutionary algorithm**

Scaling:

- Consider larger networks (10 species)