# Classical Improvements to Modern Machine Learning

Shiva Kaul <skkaul@cs.cmu.edu>

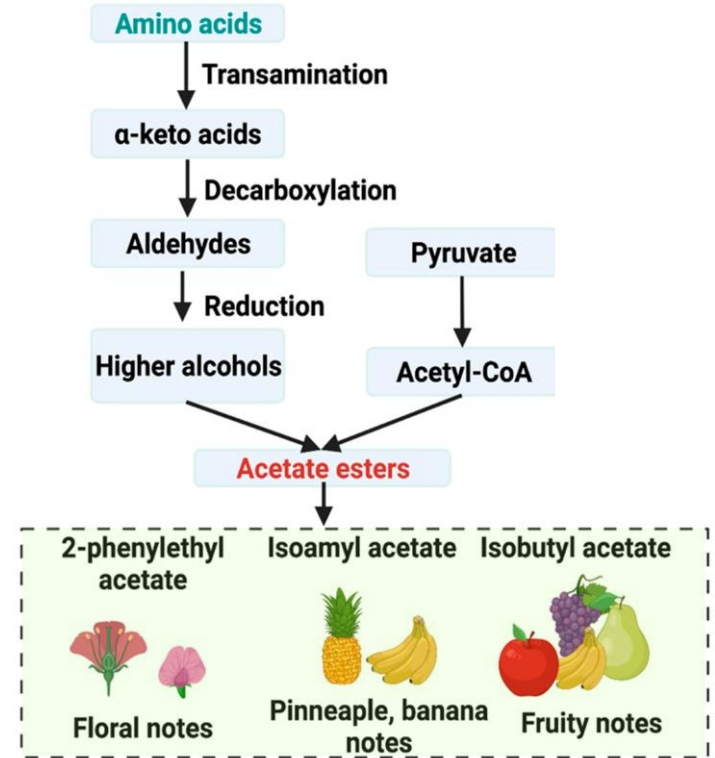Sugar metabolism

*Gutiérrez-Ríos et al. (2022)*

**FERMENTED FOODS**

**WIN-WIN**

**Inhibit pathogens, add flavor compounds**

Protein metabolism

# Utilitarian

- Nutritious
- Long-lasting
- Easy to prepare

# Hedonic
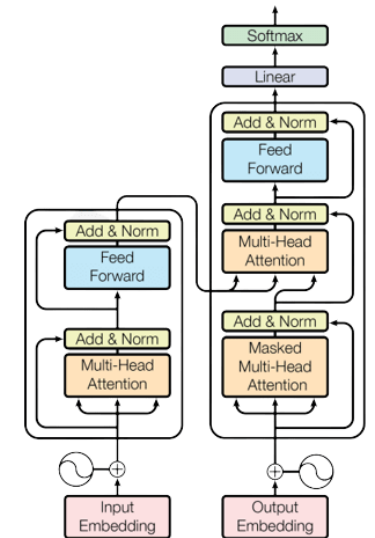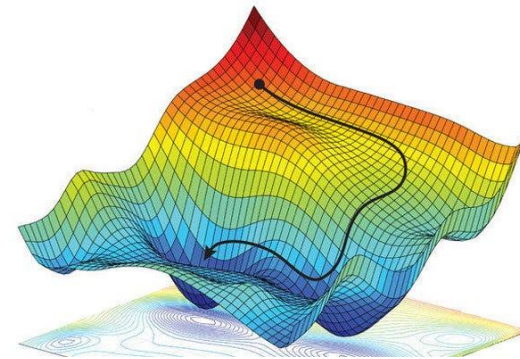
Tastes good •

# Classical

- Efficient
- Safe (reliable, robust, interpretable)
- Easy to analyze

# Modern

Accurate •

# Classical

# Modern

- Efficient
- Safe (reliable
- Easy to analy

Accurate •

# Syntheses between classical and modern machine learning

*What are the* **WIN-WIN** *mechanisms?*

# Classical

# Modern

- Efficient
- Safe (reliable
- Easy to analy

Accurate •



**Meta-Analysis**

$$f\left(\text{What is the most famous cheese in France?}\right) = \text{It is arguably Camembert.}$$

**Sequence Models**

*What are the* **WIN-WIN** *mechanisms?*

# **Meta-analysis** is an interesting machine learning problem.

- Use large datasets to dramatically improve causal inferences and patient outcomes

- Scientific question-answering is an interesting unsolved problem for LLMs

- A beautiful statistical problem which exposes key challenges in uncertainty quantification

## Meta-Analysis

Effect

$$U_i = \text{ATE} + N(0, \nu)$$

between-trial
heterogeneity

$$Y_i = U_i + N(0, V_i)$$

Observed effect

within-trial variance

ATE

$V_7$

$U_7$  $Y_7$

$v$  $y$ $u$

**Open Access**                    **Research**

**BMJ Open** Plea for routinely presenting prediction
intervals in meta-analysis

Joanna IntHout,[1] John P A Ioannidis,[2,3,4,5] Maroeska M Rovers,[1] Jelle J Goeman[1]

*Letelier et al. (2003)*

Galperin et al[29] (2000) 33.7 (2.08-546.00) 95
Bianconi et al[28] (2000) 2.04 (0.19-22.00) 83
Villani et al[11] (2000) 4.75 (1.60-14.00) 120
Hohnloser et al[3] (2000) 3.13 (1.5-6.70) 203
Natale et al[25] (2000) 5.12 (2.60-10.00) 85
Cowan et al[16] (1986) 1.11 (0.78-1.58) 34
Noc et al[17] (1990) 18.00 (1.17-276.00) 24
Capucci et al[18] (1992) 0.77 (0.37-1.62) 40
Cochrane et al[19] (1994) 1.15 (0.91-1.44) 30
Hou et al[21] (1995) 1.29 (0.97-1.72) 39
Kondili et al[22] (1995) 1.33 (0.71-2.47) 42
Donovan et al[20] (1994) 1.05 (0.69-1.60) 64
Galve et al[23] (1996) 1.13 (0.84-1.52) 100
Kontoyannis et al[24] (1998) 1.42 (1.08-1.85) 42
Bellandi et al[26] (1999) 1.41 (1.15-1.72) 120
Kochiadakis et al[12] (1999) 1.46 (1.19-1.78) 204
Cotter et al[27] (1999) 1.43 (1.15-1.8) 100
Peuhkurinen et al[30] (2000) 2.45 (1.49-4.02) 62
Vardas et al[31] (2000) 2.01 (1.55-2.6) 208
Joseph and Ward[32] (2000) 1.32 (095-1.80) 75
Cybulski et al[33] (2001) 1.87 (1.37-2.55) 160

[Future]

**95% CI**

**95% PI for y**

**95% PI for u**

0.1        1        10
Relative Risk (95% CI)

**Meta-Analysis**

Features   Effect   Variance

$$(X_i, \ U_i, \ V_i) \ \sim \ \mathbb{P}$$

$$Y_i = U_i + N(0, V_i)$$

Observed effect

$$1 - \alpha \leq \mathbb{P}(y \in C(x, v))$$

$$1 - \alpha \leq \mathbb{P}(u \in C(x))$$

$X_1$

$X_7$

Galperin et al[29] (2000) 33.7 (2.08-546.00) 95
Bianconi et al[28] (2000) 2.04 (0.19-22.00) 83
Villani et al[11] (2000) 4.75 (1.60-14.00) 120
Hohnloser et al[3] (2000) 3.13 (1.5-6.70) 203
Natale et al[25] (2000) 5.12 (2.60-10.00) 85
Cowan et al[16] (1986) 1.11 (0.78-1.58) 34
Noc et al[17] (1990) 18.00 (1.17-276.00) 24
Capucci et al[18] (1992) 0.77 (0.37-1.62) 40
Cochrane et al[19] (1994) 1.15 (0.91-1.44) 30
Hou et al[21] (1995) 1.29 (0.97-1.72) 39
Kondili et al[22] (1995) 1.33 (0.71-2.47) 42
Donovan et al[20] (1994) 1.05 (0.69-1.60) 64
Galve et al[23] (1996) 1.13 (0.84-1.52) 100
Kontoyannis et al[24] (1998) 1.42 (1.08-1.85) 42
Bellandi et al[26] (1999) 1.41 (1.15-1.72) 120
Kochiadakis et al[12] (1999) 1.46 (1.19-1.78) 204
Cotter et al[27] (1999) 1.43 (1.15-1.8) 100
Peuhkurinen et al[30] (2000) 2.45 (1.49-4.02) 62
Vardas et al[31] (2000) 2.01 (1.55-2.6) 208
Joseph and Ward[32] (2000) 1.32 (095-1.80) 75
Cybulski et al[33] (2001) 1.87 (1.37-2.55) 160

**1-α PI for y**

**1-α PI for u**

0.1   1   10
Relative Risk (95% CI)

# Trusted data

- Rigorous and unbiased
- Loose predictions



# Untrusted data

Need strong assumptions •

Tight predictions •

# Trusted data

- Rigorous and unbiased
- Loose predictions

# Untrusted data

- Need strong assumptions
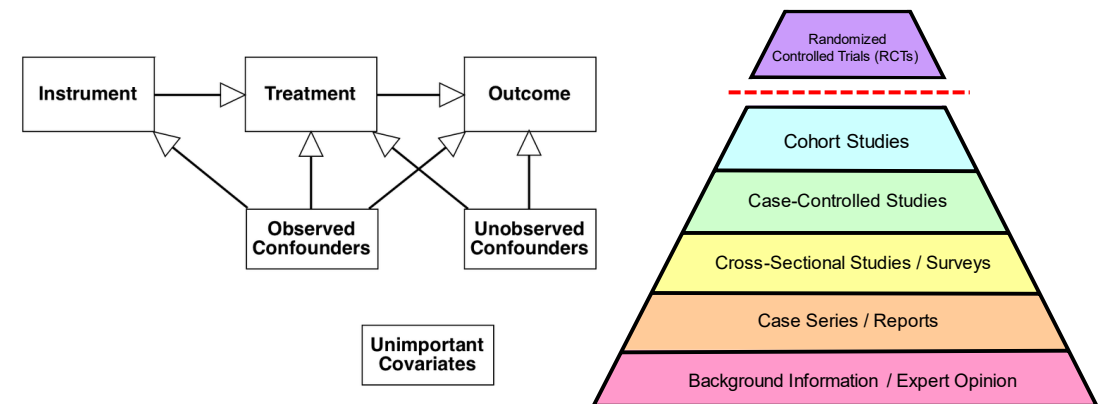- Tight predictions

## CONFORMAL META-ANALYSIS



Systematic Reviews

$C$

**CONFORMAL META-ANALYSIS**

Randomized Controlled Trials (RCTs)

trusted data

$X_i, Y_i, V_i$

Cohort Studies

untrusted data

$\mu, \kappa$

Case-Controlled Studies

Cross-Sectional Studies / Surveys

Case Series / Reports

Background Information / Expert Opinion

- - - - - - - - - - - - - - - WIN-WIN - - - - - - - - - - - - - - -

**Assumption-free inclusion of untrusted prior**

$$0.90 \leq \mathbb{P}(r^* \text{ among lowest 20 of } R_i)$$

$$C(x, v) = \{y : r \text{ among lowest 20 of } R_i\}$$

$\longrightarrow$

$$0.90 \leq \mathbb{P}(y^* \in C(x, v))$$

**Conformal Prediction**

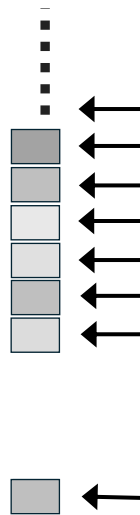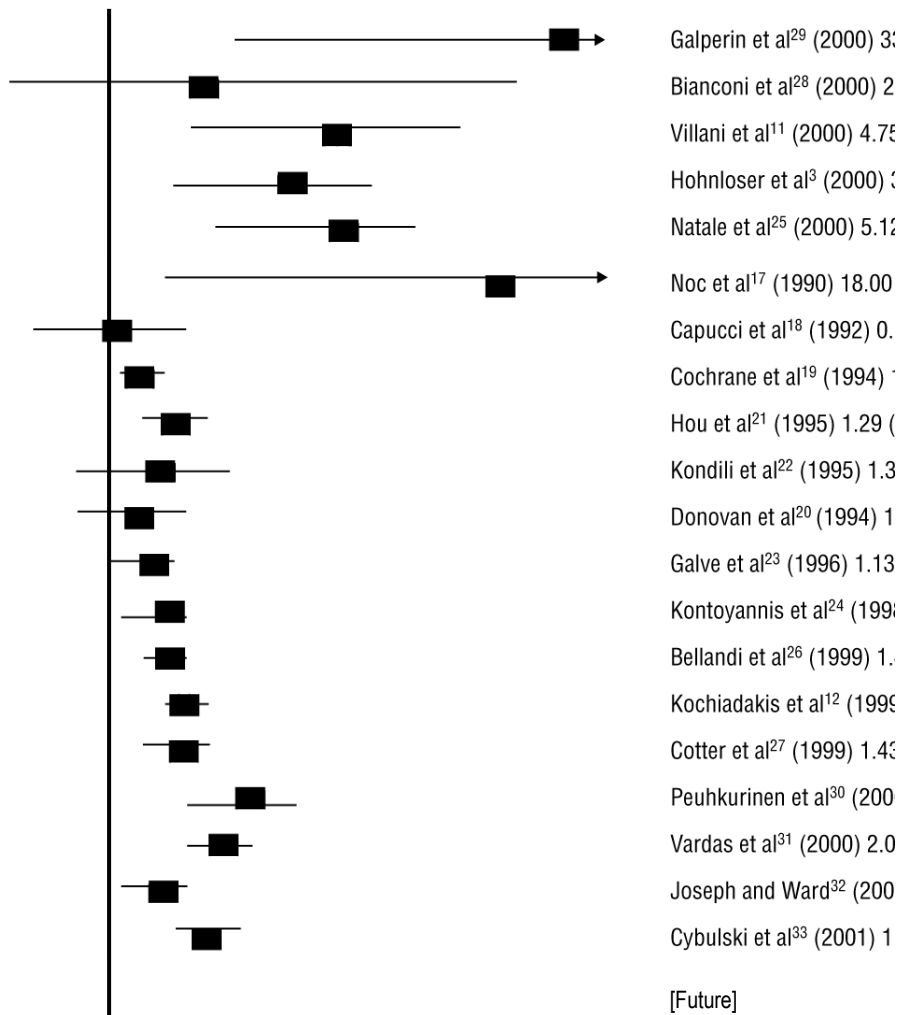$$n = 21$$

**Past** $(X_i, Y_i, V_i)$ has residual $R_i$

*exchangeable, so rank of r\* is uniform among $R_i$*

**Future** $(x, y^*, v)$ has residual $r^*$

**(Hypothetical) Future** $(x, y, v)$ has residual $r$

Galperin et al[29] (2000) 3:

Bianconi et al[28] (2000) 2

Villani et al[11] (2000) 4.75

Hohnloser et al[3] (2000) :

Natale et al[25] (2000) 5.12

Noc et al[17] (1990) 18.00

Capucci et al[18] (1992) 0.

Cochrane et al[19] (1994)

Hou et al[21] (1995) 1.29 (

Kondili et al[22] (1995) 1.3

Donovan et al[20] (1994) 1

Galve et al[23] (1996) 1.13

Kontoyannis et al[24] (199:

Bellandi et al[26] (1999) 1.

Kochiadakis et al[12] (1999

Cotter et al[27] (1999) 1.4:

Peuhkurinen et al[30] (200

Vardas et al[31] (2000) 2.0

Joseph and Ward[32] (200

Cybulski et al[33] (2001) 1

[Future]

**90% PI for y**      ??????????????????      $C(x, v) = \{y : r \text{ among lowest 20 of } R_i\}$

**Prior** $\mu, \kappa$ ——$[Y; y]$—$[V; v]$——$[X; x]$→ **Posterior** $\hat{\mu}, \hat{\kappa}$ → **Residuals** $R_i, r$



Galperin et al[29] (2000) 3?

Bianconi et al[28] (2000) 2

Villani et al[11] (2000) 4.75

Hohnloser et al[3] (2000) ?

Natale et al[25] (2000) 5.1?

Noc et al[17] (1990) 18.00

Capucci et al[18] (1992) 0.

Cochrane et al[19] (1994)

Hou et al[21] (1995) 1.29 (

Kondili et al[22] (1995) 1.3

Donovan et al[20] (1994) 1

Galve et al[23] (1996) 1.13

Kontoyannis et al[24] (199?

Bellandi et al[26] (1999) 1.

Kochiadakis et al[12] (199?

Cotter et al[27] (1999) 1.4?

Peuhkurinen et al[30] (200

Vardas et al[31] (2000) 2.0

Joseph and Ward[32] (200

Cybulski et al[33] (2001) 1

[Future]

**90% PI for y**   ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ? ?    $C(x, v) = \{y : r \text{ among lowest 20 of } R_i\}$

$y$
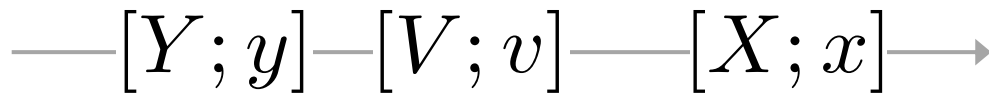
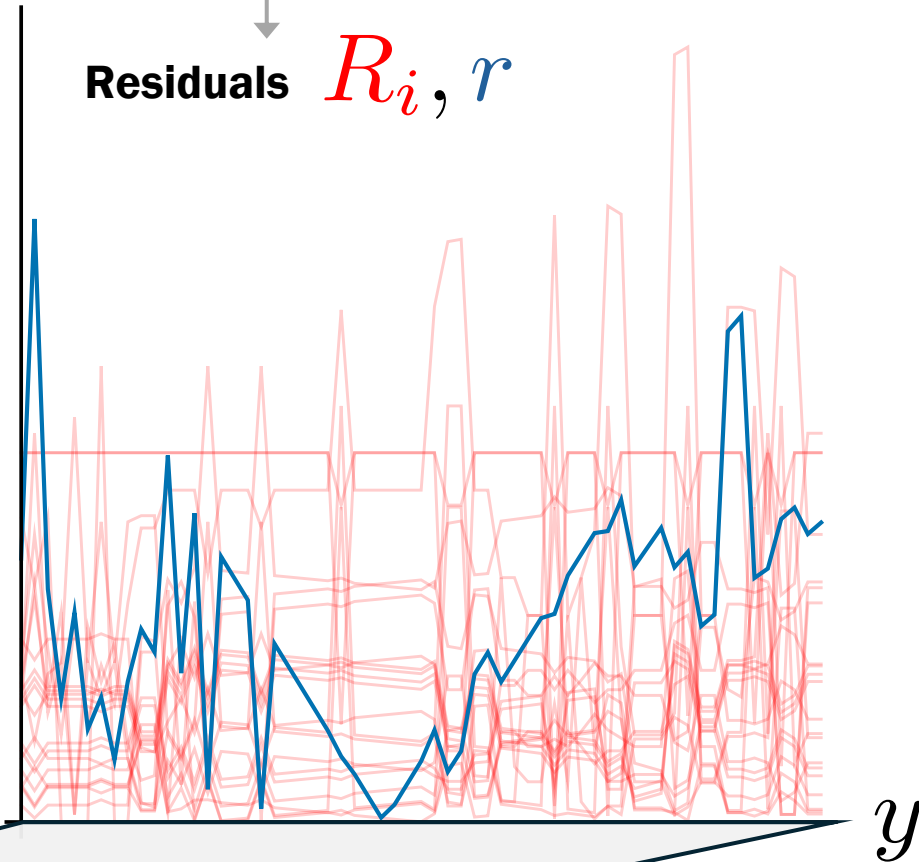train on everything for exchangeability

**Prior** $\mu, \kappa$ —— $[Y; y]$—$[V; v]$——$[X; x]$ ⟶ **Posterior** $\hat{\mu}, \hat{\kappa}$ ⟶ **Residuals** $R_i, r$

Galperin et al[29] (2000) 3:
Bianconi et al[28] (2000) 2
Villani et al[11] (2000) 4.75
Hohnloser et al[3] (2000) 3
Natale et al[25] (2000) 5.12
Noc et al[17] (1990) 18.00
Capucci et al[18] (1992) 0.
Cochrane et al[19] (1994)
Hou et al[21] (1995) 1.29 (
Kondili et al[22] (1995) 1.3
Donovan et al[20] (1994) 1
Galve et al[23] (1996) 1.13
Kontoyannis et al[24] (199
Bellandi et al[26] (1999) 1.
Kochiadakis et al[12] (1999
Cotter et al[27] (1999) 1.43
Peuhkurinen et al[30] (200
Vardas et al[31] (2000) 2.0
Joseph and Ward[32] (200
Cybulski et al[33] (2001) 1
[Future]

$y$

**90% PI for y**  ??????????????????  $C(x, v) = \{y : r \text{ among lowest 20 of } R_i\}$

Prior $\mu, \kappa$

*train on everything for exchangeability*

$[Y; y] - [V; v] - [X; x] \longrightarrow$ Posterior $\hat{\mu}, \hat{\kappa}$

Residuals $R_i, r$

Galperin et al[29] (2000) 3

Bianconi et al[28] (2000) 2

Villani et al[11] (2000) 4.75

Hohnloser et al[3] (2000)

Natale et al[25] (2000) 5.12

Noc et al[17] (1990) 18.00

Capucci et al[18] (1992) 0.

Cochrane et al[19] (1994)

Hou et al[21] (1995) 1.29 (

Kondili et al[22] (1995) 1.3

Donovan et al[20] (1994) 1

Galve et al[23] (1996) 1.13

Kontoyannis et al[24] (199

Bellandi et al[26] (1999) 1.

Kochiadakis et al[12] (1999

Cotter et al[27] (1999) 1.43

Peuhkurinen et al[30] (200

Vardas et al[31] (2000) 2.0

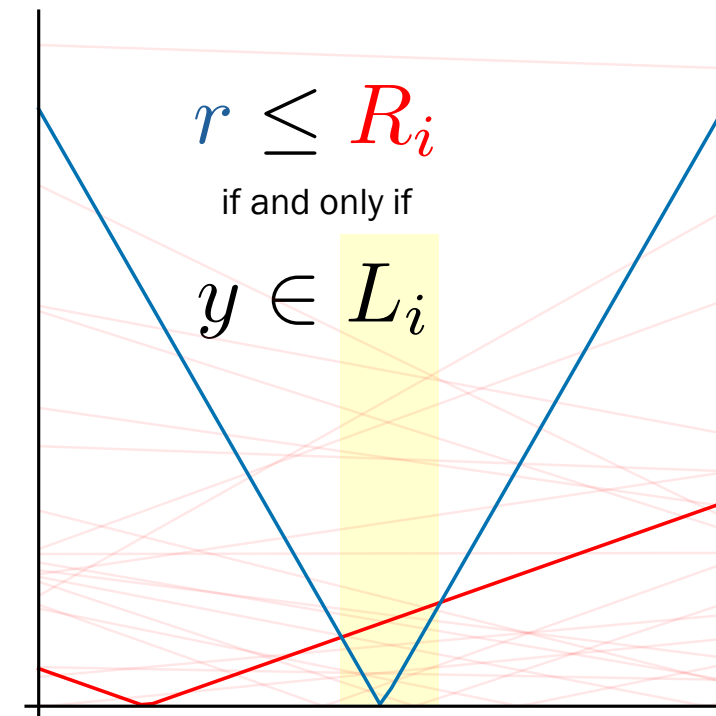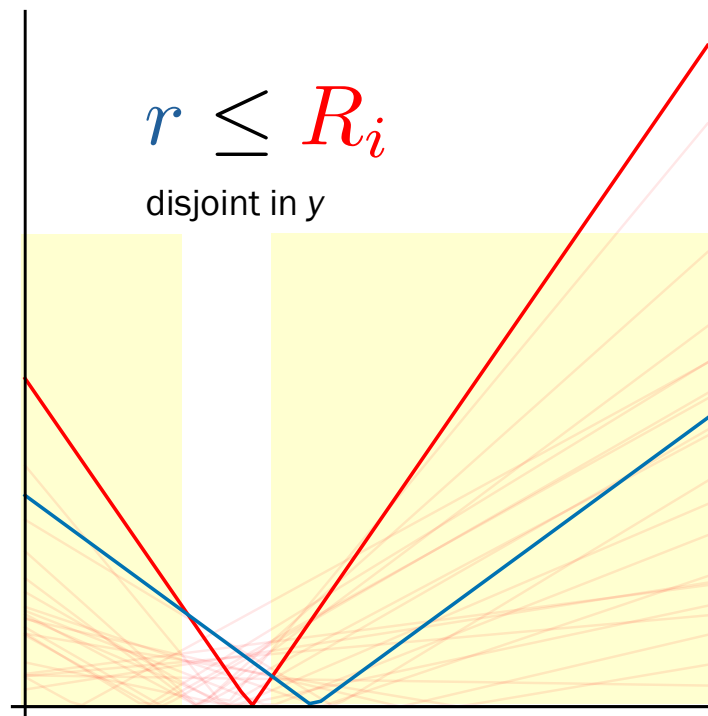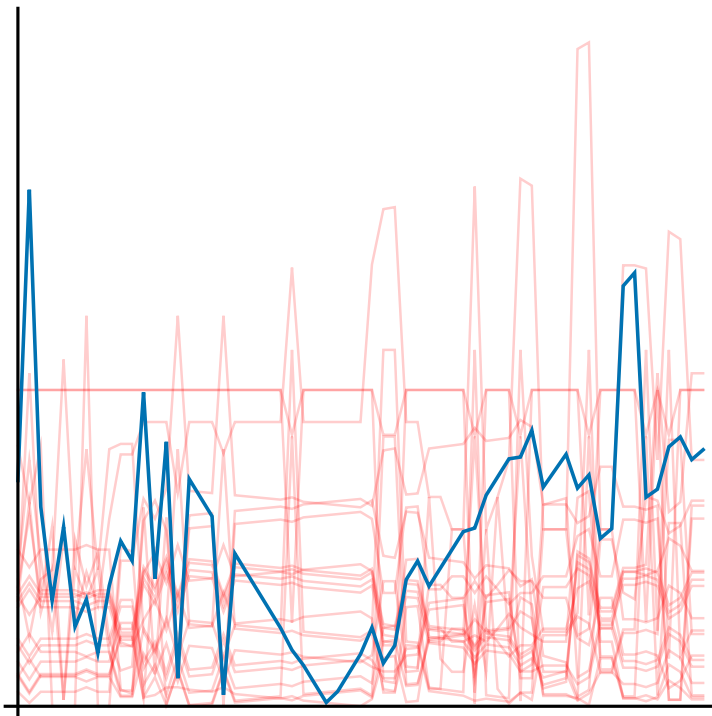Joseph and Ward[32] (20

Golluski et al[33] (2001) 1

[Future]

$y$

**90% PI for y**

$C(x, v) = \{y : r \text{ among lowest 20 of } R_i\}$

## *CHALLENGES*

**1.** Full conformal prediction is intractable

($n$ is small, so cannot split the data)

**2.** Also want interval for $u$, not just $y = N(u,v)$

*Kaul and Gordon (2024)*

Focus on **linear smoothers**

like kernel ridge regression (KRR)

$$r \leq R_i$$
disjoint in $y$

$$r \leq R_i$$
if and only if
$$y \in L_i$$

$$R_i = \ldots |A_i y + B_i| \ldots \qquad r = \ldots |ay + b| \ldots$$

residuals are convex in $y$

Ensure **idiocentricity**

changing $y$ affects $r$ more than any $R_i$

$$|a| > |A_i| \quad \Longleftarrow \quad \lambda \geq \max_x \kappa(x, x)$$

for linear smoothers                     easy to ensure for KRR

\*Tolerate **approximation**

$$C(x, v) \subseteq \left[ \begin{array}{cc} \text{2nd lowest} & \text{2nd highest} \\ \text{left end of } L_i & , \quad \text{right end of } L_i \end{array} \right]$$
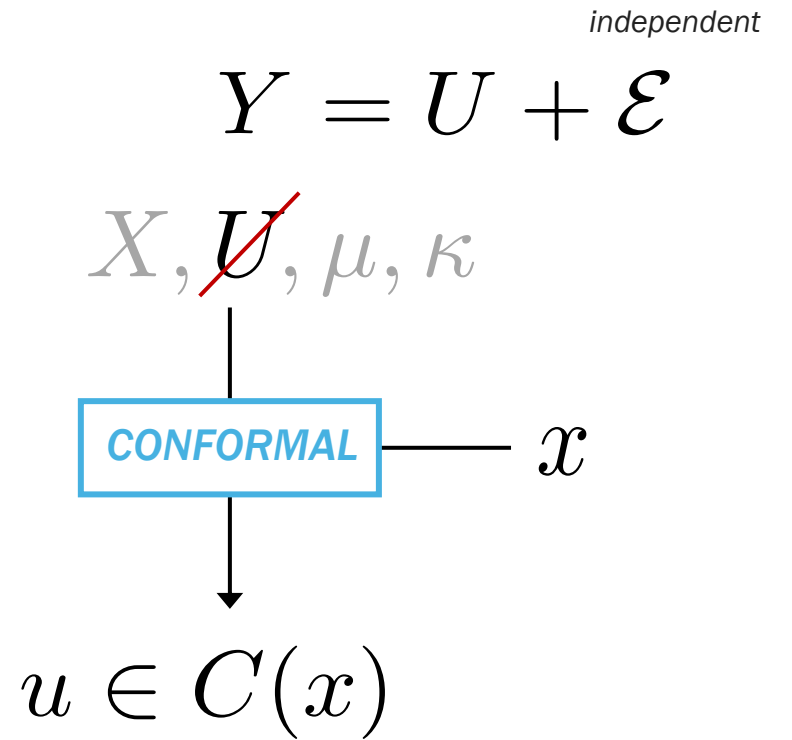
## *CHALLENGES*

**1.** Full conformal prediction is intractable

**...**but not for idiocentric linear smoothers.

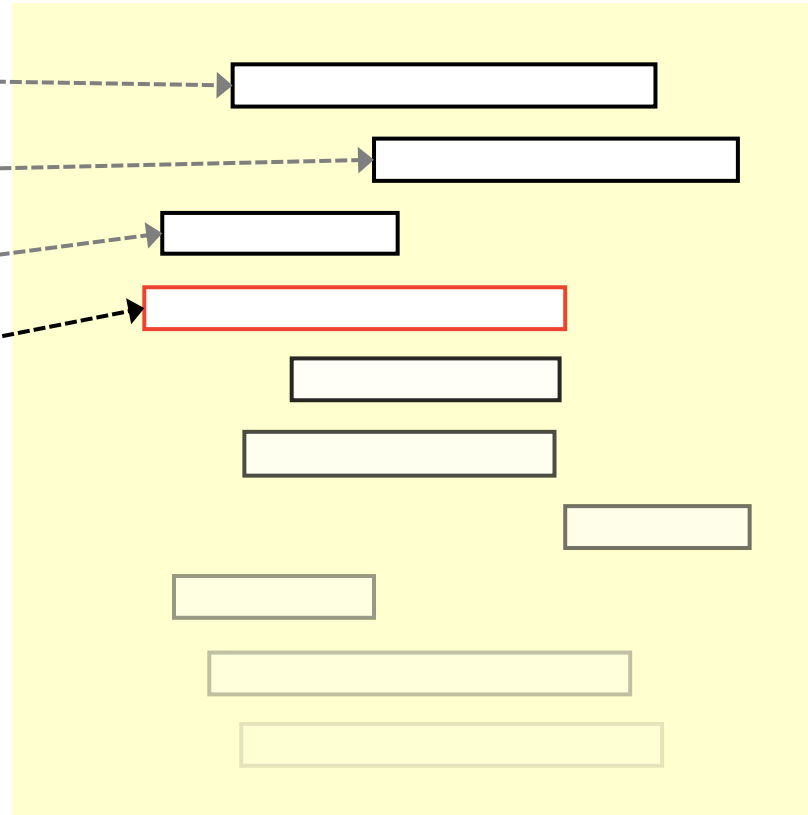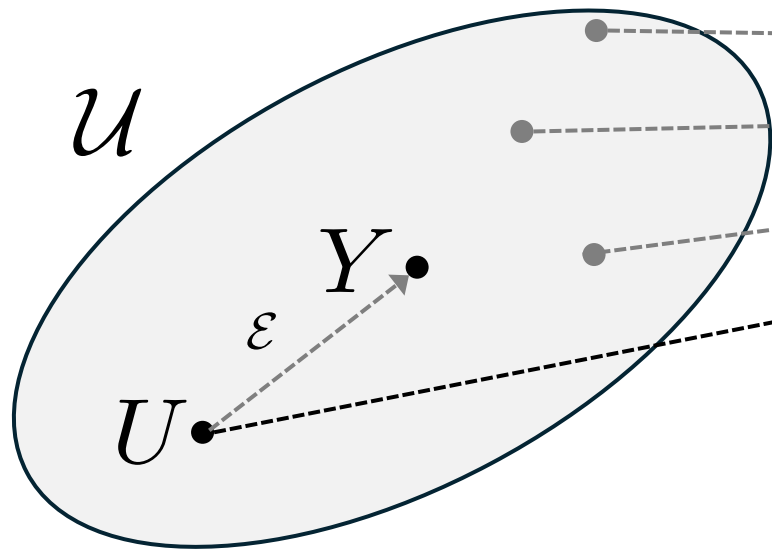**2.** Also want interval for *u*, not just *y* = *N(u,v)*

*Kaul and Gordon (2024)*

$$Y = U + \mathcal{E}$$

$$\mu, \kappa, X, Y, V$$

$$X, \cancel{U}, \mu, \kappa$$

| CONFORMAL |

$x, v$ —

— $x$

$$y \in C(x, v)$$

**VS**

$$u \in C(x)$$

$$\mathbb{P}_{\mathcal{E}}(U \in \mathcal{U}) \geq 1 - \delta \qquad \bigcup \{C(x; \widehat{U}) : \widehat{U} \in \mathcal{U}\}$$
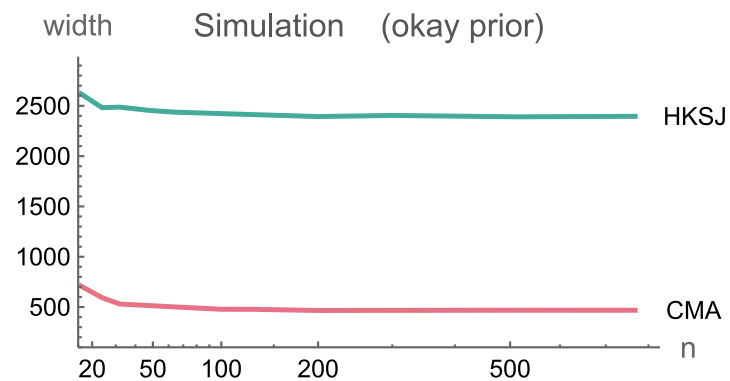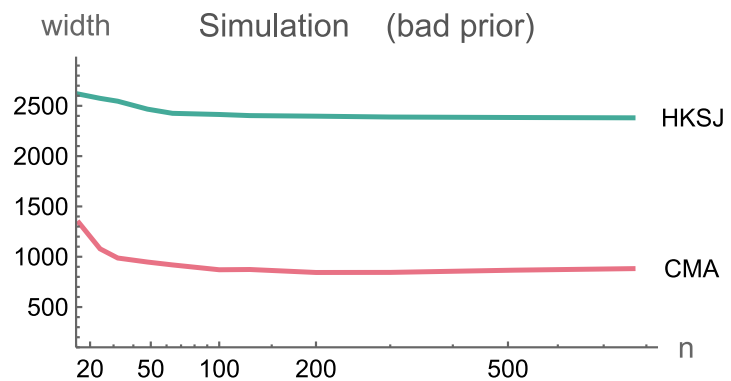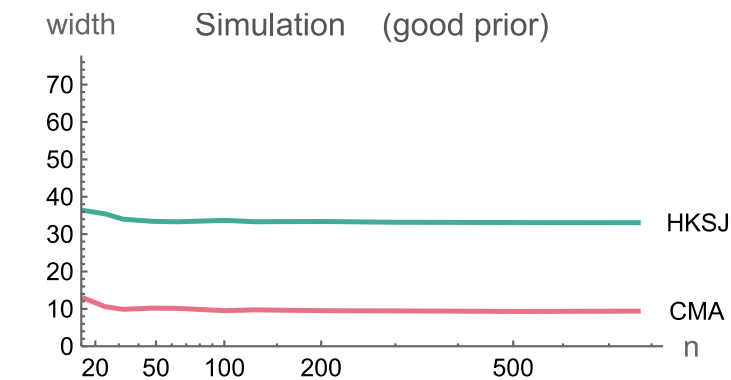


$\mathcal{U}$

$Y$

$\mathcal{E}$

$U$

$$(1 - \alpha)(1 - \delta) \leq \mathbb{P}\left(u \in [ \qquad \qquad ]\right)$$

Exploit independence of noise $\mathcal{E}$

Idiocentricity → tightly bound outer interval

- Formulated meta-analysis as an interesting machine learning problem

- Simplified full conformal prediction for idiocentric linear smoothers

- Addressed statistical/algorithmic challenges in handling noise

SPECIAL ARTICLE

Problems in the ''Evidence'' of ''Evidence-based Medici

Alvan R. Feinstein, MD, Ralph I. Horwitz, MD, *New Haven, Connecticut*

The proposed practice of "evidence-based medicine," which calls for careful clinical judgment in evaluating the "best available evidence," should be differentiated from the special collection of data regarded as suitable evidence. Although the proposed practice does not seem new, the new collection of "best available" information has major constraints for

Within 5 yea based med astic endorseme cal journals,[2] acl own new journa often accorded t

Hardly anyone ting clinicians to

Above averaging in literature reviews

Uri Simonsohn [1 ✉], Joseph Simmons[2] and Leif D. Nelson[3]

Meta-analysts' practice of transcribing and numerically combining all results in a research literature can generate uninterpretable and/or misleading conclusions. Meta-analysts should

ovide

Perspective | Published: 22 July 2021

Behavioural science is unlikely to change the world without a heterogeneity revolution

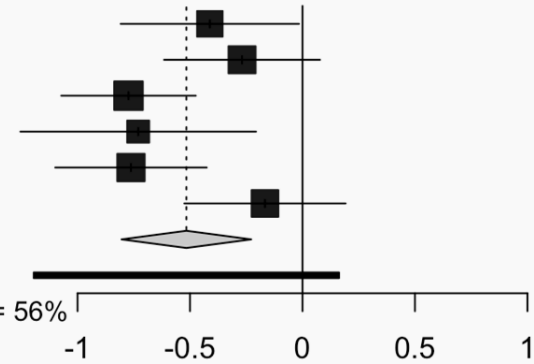Christopher J. Bryan ✉ Elizabeth Tipton ✉ & David S. Yeager ✉

iefly men-lance as to , such rec-ce, studies which are or studies

**Meta-Analysis**

$$f\left(\text{What is the most famous cheese in France?}\right) = \text{It is arguably Camembert.}$$

**Sequence Models**

# Linear

- Efficient (*O(T)* memory) and
- Fast (*O(log T)* parallel time via scans)
- Unexpressive

$$h_t = A_t h_{t-1} + B_t x_t$$

*(Time-varying) Linear dynamical system*

# Nonlinear

Inefficient (*O(T²)* memory) or •
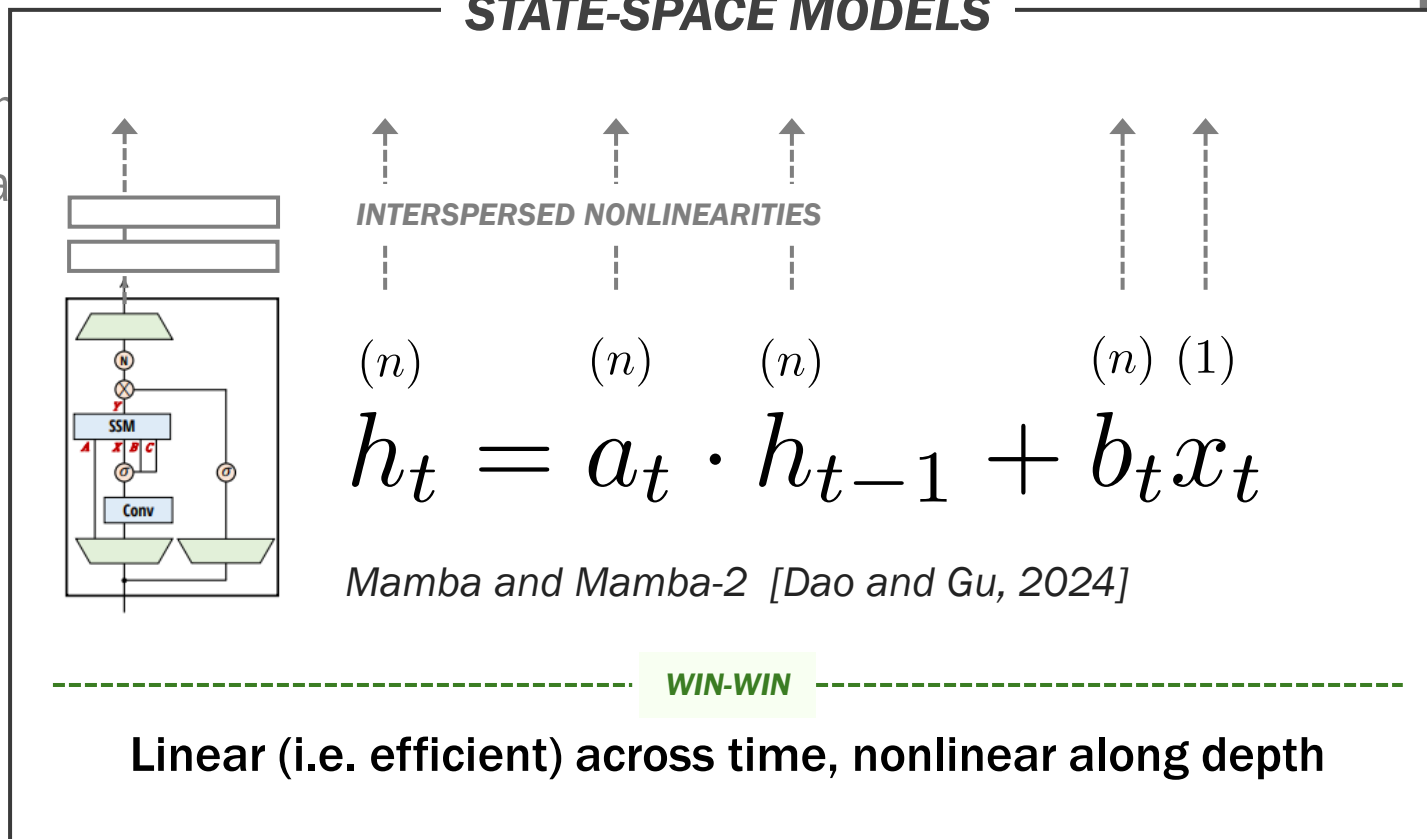Slow (*O(T)* parallel time) •
Expressive •

$$h = \psi(Q(x)K(x))V(x)$$

*Attention*

$$\overset{(n)}{h_t} = \rho(\overset{(n,n)}{A_t} \overset{(n)}{h_{t-1}} + \overset{(n,d)}{B_t} \overset{(d)}{x_t})$$

*Recurrent neural network*

# Linear

- Efficient (*O(T)* mem
- Fast (*O(log T)* para
- Unexpressive

# Nonlinear

- (*O(T²)* memory) or
- *O(T)* parallel time)
- Expressive

## STATE-SPACE MODELS



*INTERSPERSED NONLINEARITIES*

$$h_t \overset{(n)}{=} a_t \overset{(n)}{\cdot} h_{t-1} \overset{(n)}{+} b_t \overset{(n)}{x_t} \overset{(1)}{}$$

$$h_t = \overset{(n)}{a_t} \cdot \overset{(n)}{h_{t-1}} + \overset{(n)}{b_t} \overset{(1)}{x_t}$$

*Mamba and Mamba-2  [Dao and Gu, 2024]*

**WIN-WIN**

**Linear (i.e. efficient) across time, nonlinear along depth**

$$K(x))V(x)$$

$$h_t = A_t h_{t-1} + B_t x_t$$

*(Time-varying) Linear dynamical system*

$$h_t = \overset{(n)}{\rho}(\overset{(n,n)}{A_t} \overset{(n)}{h_{t-1}} + \overset{(n,d)}{B_t} \overset{(d)}{x_t})$$

*Recurrent neural network*

# Nonlinearity ~~across time~~ along depth via iterated local corrections *[Kaul 2020]*

*Goal: approximate nonlinear RNN by a stack of linear systems, with nonlinearity along only depth*

**Theory:** understand power of depth

**Practice:** use within new models

## The Illusion of State in State-Space Models

## Theoretical Foundations of Deep Selective State-Space Models

Nicola Muca Cirone [1]   Antonio Orvieto [2]   Benjamin Walker [3]   Cristopher Salvi [1]   Terry Lyons [3]

### Abstract

Structured state-space models (SSMs) such as S4, stemming from the seminal work of Gu et al., are gaining popularity as effective approaches for modeling sequential data. Deep SSMs demonstrate outstanding performance across a diverse achieve state-of-the-art results on long-range-reasoning benchmarks (Tay et al., 2020) and show outstanding performance in various domain including vision (Nguyen et al., 2022), audio (Goel et al., 2022), biological signals (Gu et al., 2021), reinforcement learning (Lu et al., 2023) and online learning (Zucchet et al., 2023). SSMs recently have gained

## Mamba: Linear-Time Sequence Modeling with Selective State Spaces

## Transformers are SSMs: Generalized Models and Efficient Algorithms Through Structured State Space Duality
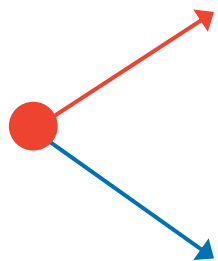
Tri Dao [*1] and Albert Gu [*2]

[1] Department of Computer Science, Princeton University

[2] Machine Learning Department, Carnegie Mellon University
tri@tridao.me, agu@cs.cmu.edu

# Nonlinearity ~~across time~~ along depth via iterated local corrections [Kaul 2020]

$$s_0^{(1)} = s_0 = h_0$$

*If a state is correct...*

$$h_1 = \rho(a_1 \cdot h_0 + b_1 x_1)$$

$$s_1 = a_1 \cdot s_0 + b_1 x_1$$

*Then its next-state multiplier is correct...*

$$k_1 = \frac{\rho(a_1 \cdot s_0 + b_1 x_1)}{a_1 \cdot s_0 + b_1 x_1} = \frac{\rho(a_1 \cdot h_0 + b_1 x_1)}{a_1 \cdot h_0 + b_1 x_1} = \frac{h_1}{a_1 \cdot h_0 + b_1 x_1}$$

*So, in the next layer, the next state becomes correct.*

$$s_1^{(1)} = k_1 \cdot \left(a_1 \cdot s_0^{(1)} + b_1 x_1\right)$$

$$= k_1 \cdot \left(a_1 \cdot h_0 + b_1 x_1\right) = h_1$$

# Nonlinearity ~~across time~~ along depth via iterated local corrections *[Kaul 2020]*

$$s_0^{(1)} = s_0 = h_0$$

$$h_t = \rho(a_1 \cdot h_{t-1} + b_1 x_1)$$

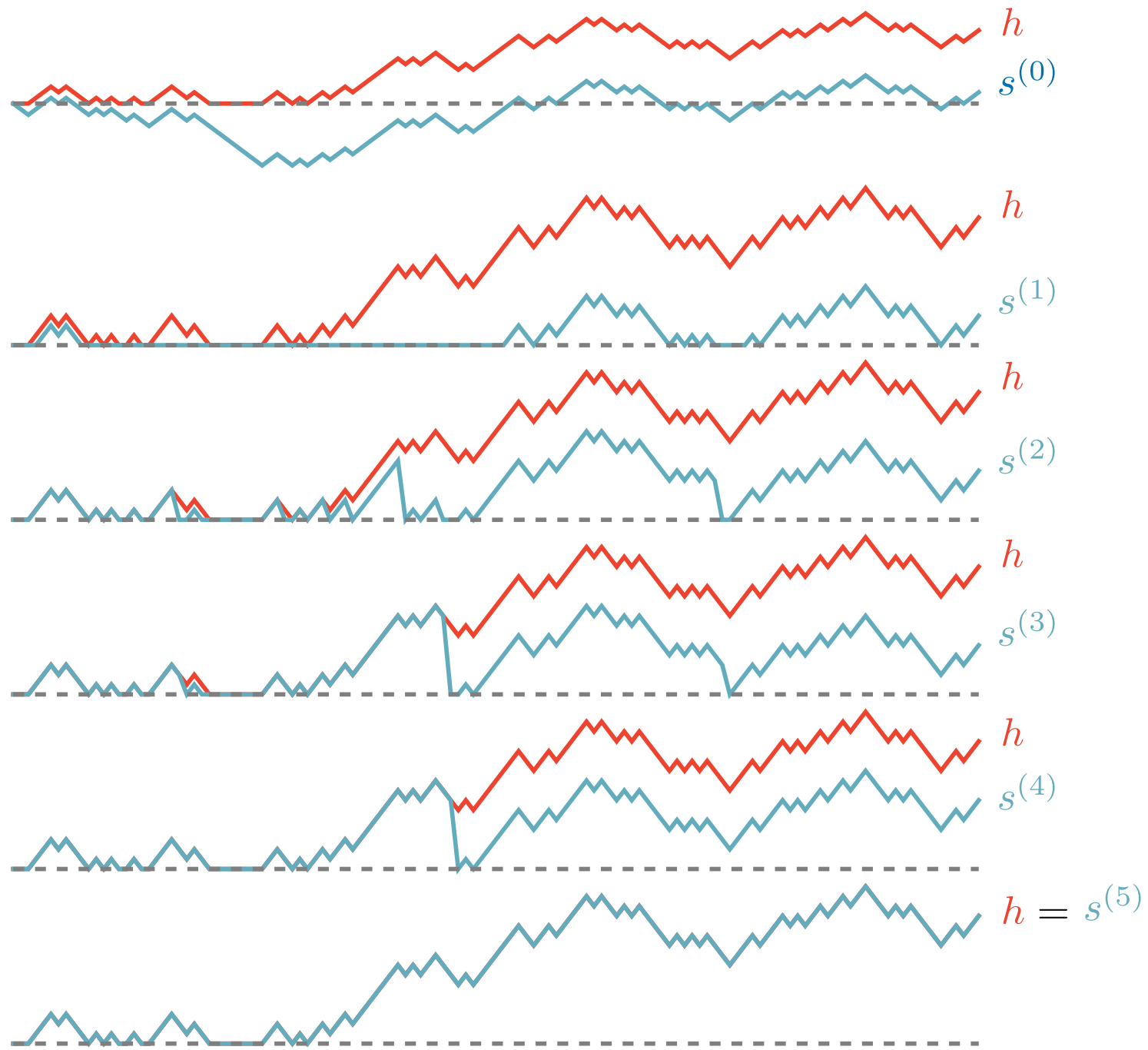$$s_t^{(0)} = a_t \cdot s_{t-1}^{(0)} + b_t x_t$$

*If a state is correct...*

$$k_t^{(i)} = \frac{\rho(a_t s_{t-1}^{(i-1)} + b_t x_t)}{a_t s_{t-1}^{(i-1)} + b_t x_t} \qquad k_i^{(i)} = \frac{h_i}{a_i h_{i-1} + b_i x_i}$$

*Then its next-state multiplier is correct...*

$$s_t^{(i)} = k_t^{(i)} \cdot (a_t \cdot s_{t-1}^{(i-1)} + b_t x_t)$$

*So, in the next layer, the next state becomes correct.*

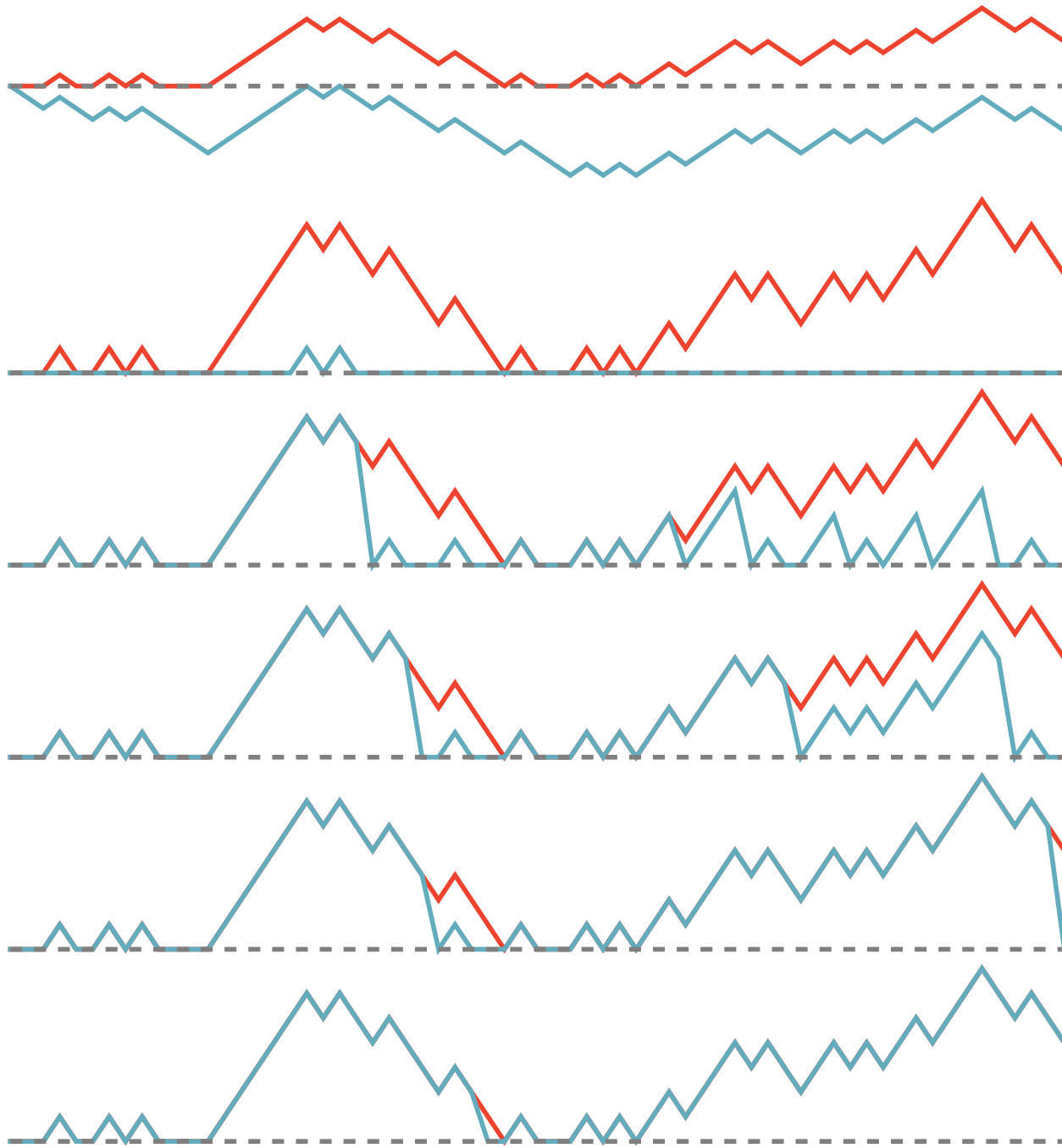$$s_i^{(i)} = k_i^{(i)} \cdot (a_i \cdot h_{i-1} + b_i x_i) = h_i$$

$\rho = \text{ReLU}$

$h$

$s^{(0)}$

$h$

$s^{(1)}$
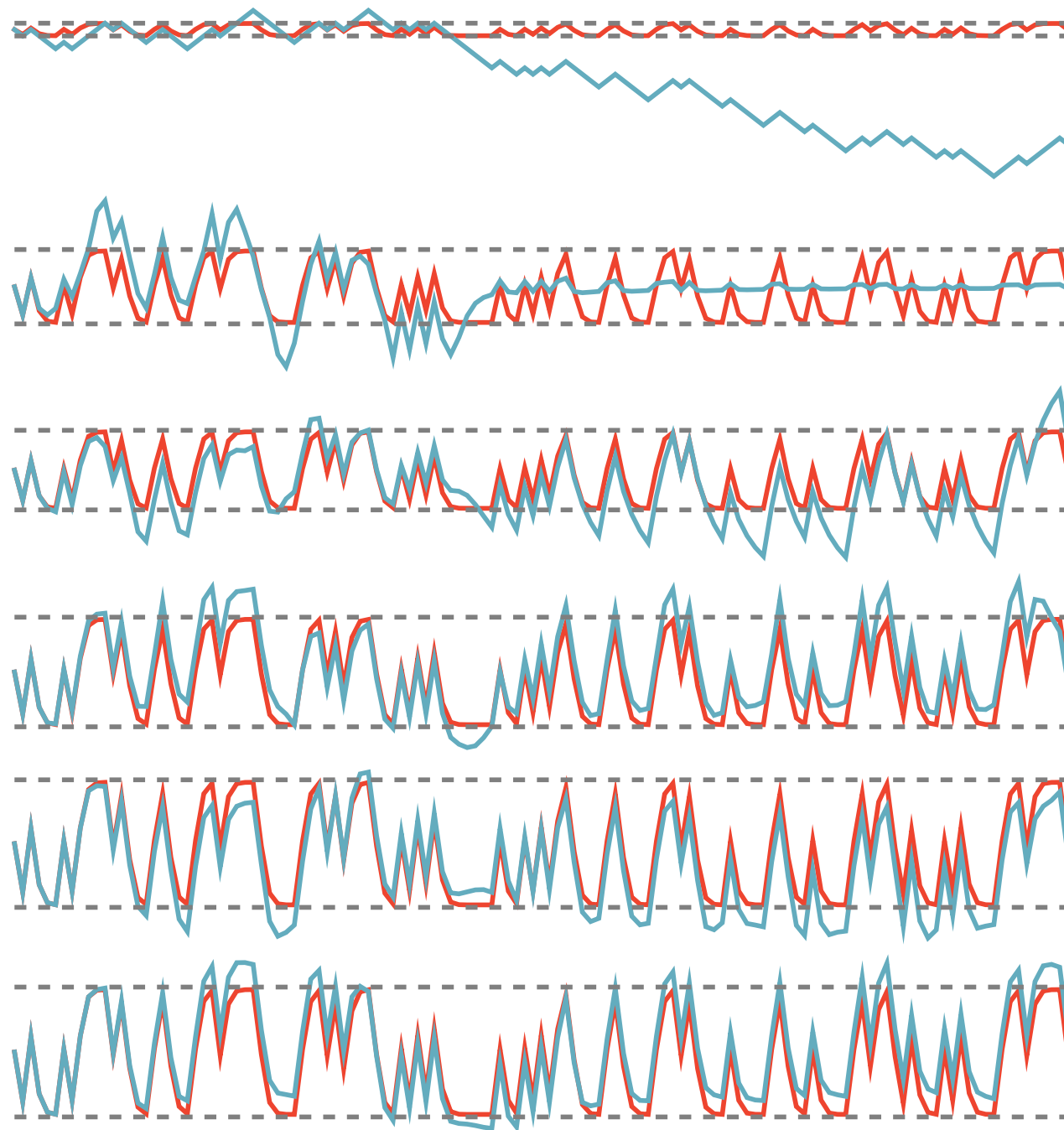
$h$

$s^{(2)}$

$h$

$s^{(3)}$

$h$

$s^{(4)}$

$h = {}^{s^{(5)}}$

$\rho = \text{ReLU}$

$\rho = \tanh$

$\rho = \text{clip}$