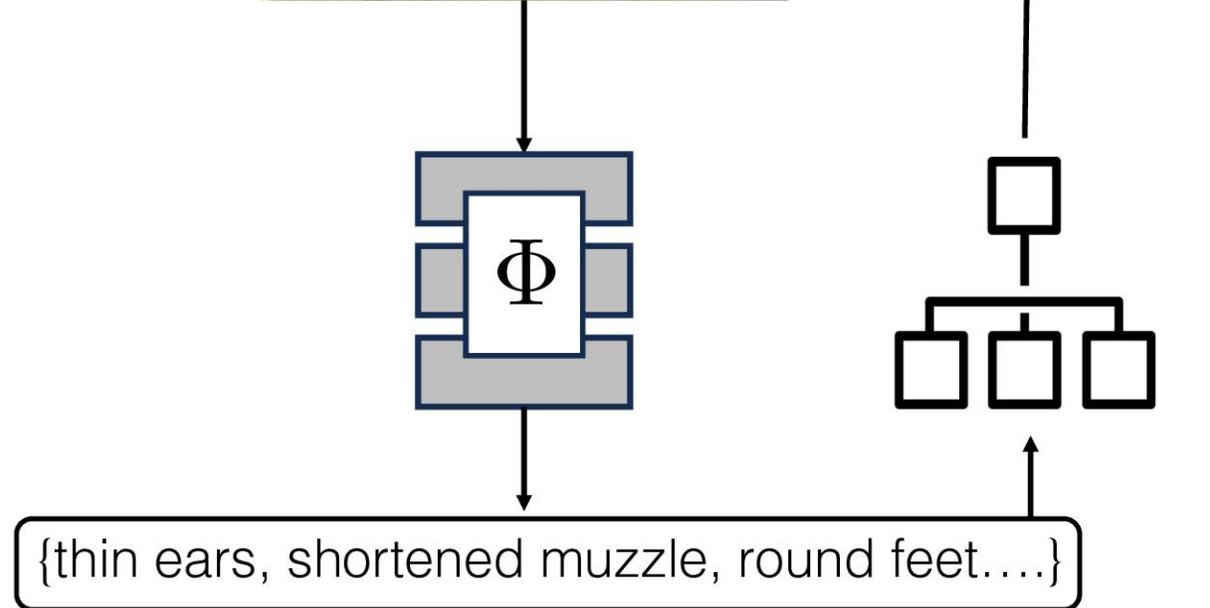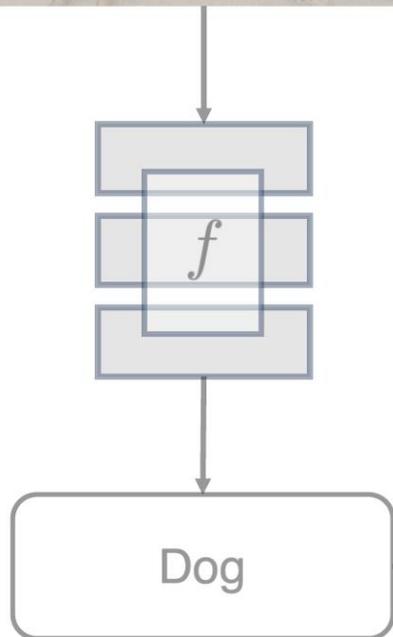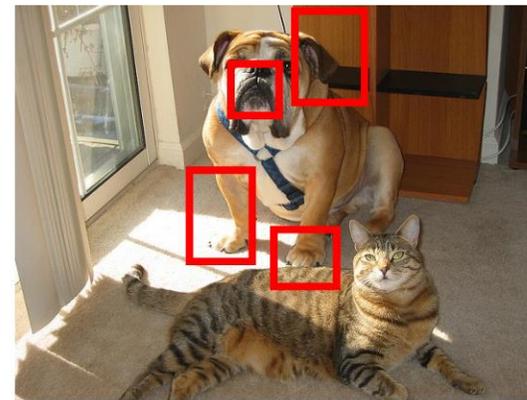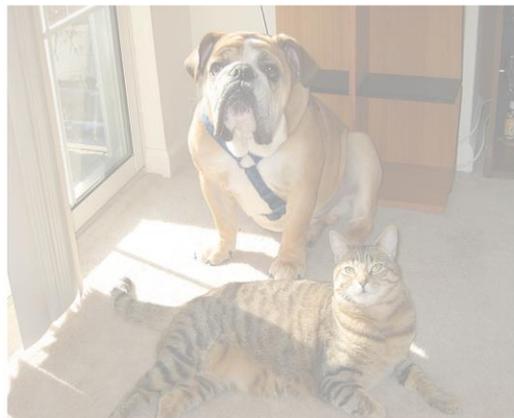# Dividing and Conquering a BlackBox to a Mixture of Interpretable Models

December 12, 2023

Seoul National University

# Post-hoc Explanations    Interpretable by design
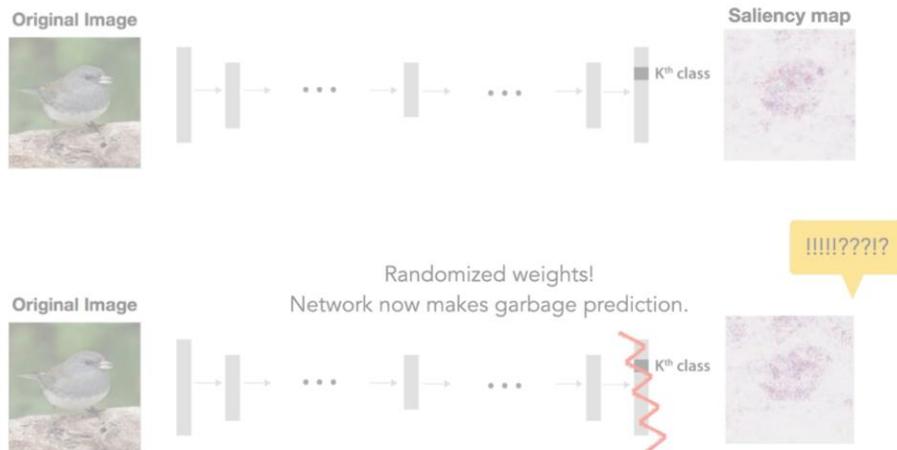


Dog

Φ

Dog

{thin ears, shortened muzzle, round feet….}

GRAD-CAM. Selvaraju et al. ICCV 2017.    Concept Bottleneck Model. Koh et al. ICML 2020

# Post-hoc explanations Interpretable by design

Pros:
- Does not alter the Black box

Cons:
- Inconsistent explanations
- No recourse

Pros:
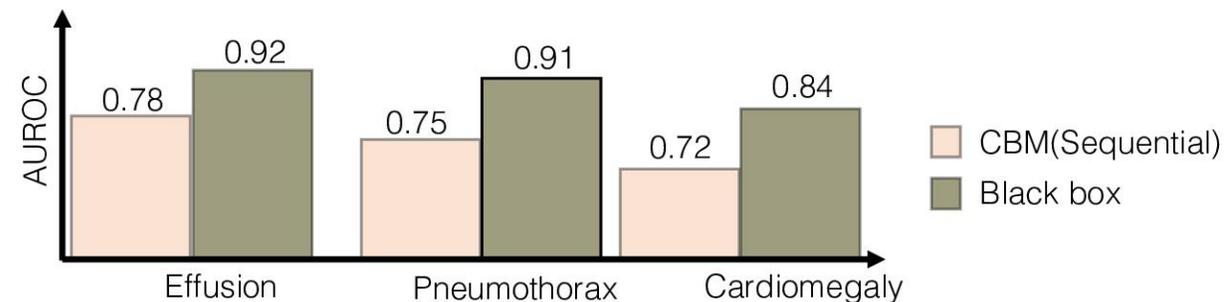- Support concept intervention

Cons:
- Harder to train
- Sub par performance

Stop explaining Black Box. Rudin et al. arXiv 2019.
Sanity Checks for Saliency Maps. Adebayo et al. Neurips 2018.

# Post-hoc explanations Interpretable by design

Pros:
• Does not alter the Black box

Pros:
• Support concept intervention

Cons:
• Inconsistent
• No recourse

Can we blur the line b/w
post-hoc explanations
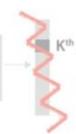or
interpretable by design

Original Image

Randomized weights!
Network now makes garbage prediction.

Original Image

!!!!!???!?

K[th] class

0.84

AUR

0.75

0.72

CBM(Sequential)
Black box

Effusion          Pneumothorax     Cardiomegaly

Stop explaining Black Box. Rudin et al. arXiv 2019.
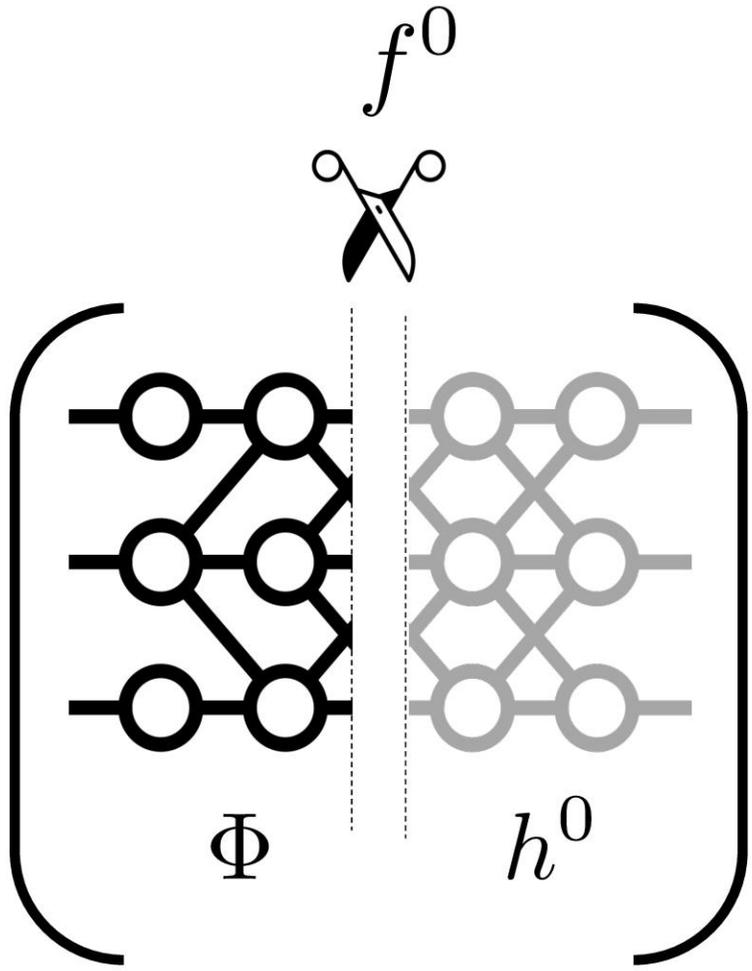Sanity Checks for Saliency Maps. Adebayo et al. Neurips 2018.

# Desirable properties

1. Does not compromise the performance
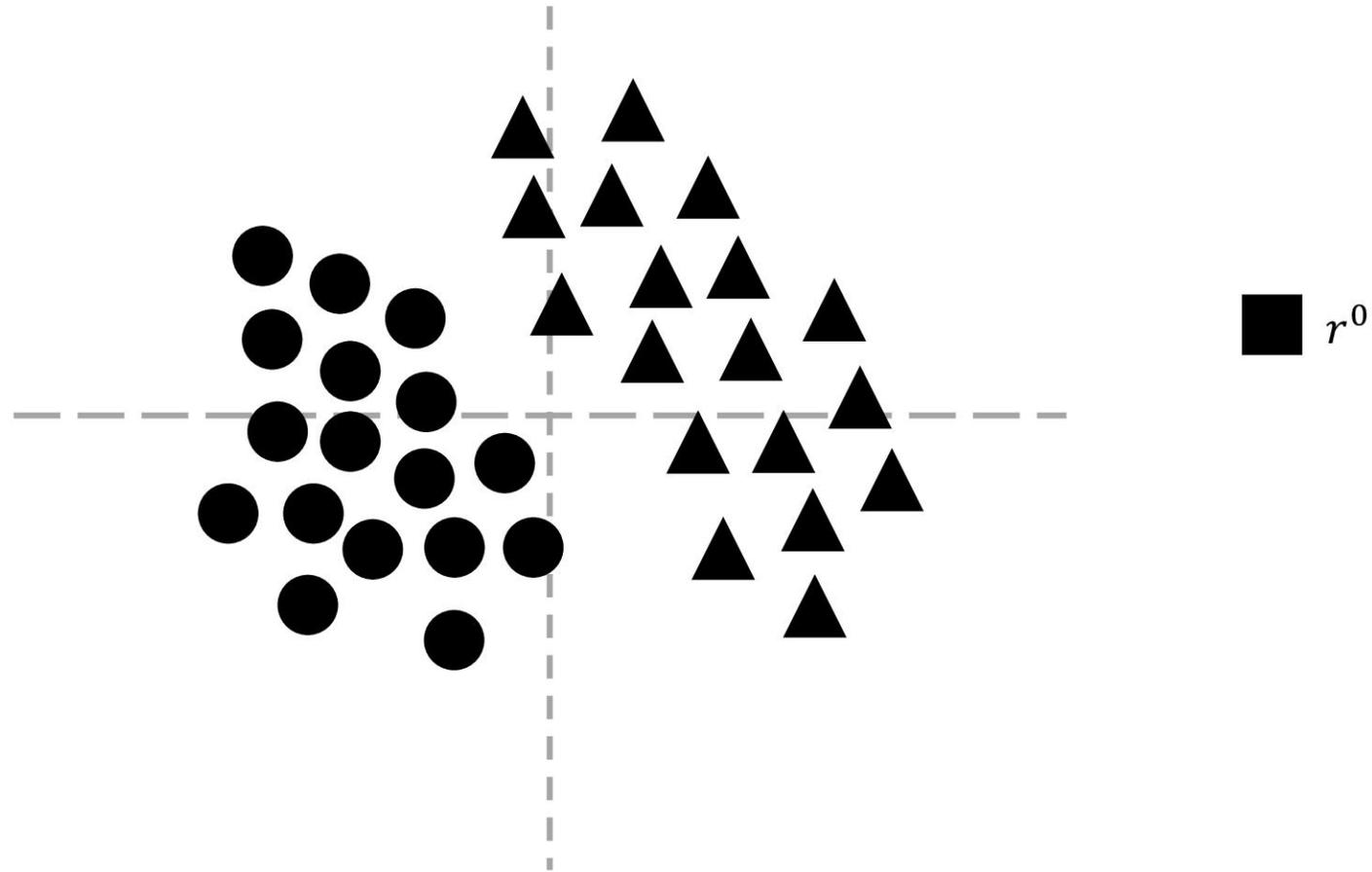2. Can be intervened to fix misclassification

# Design choices

1. Iteratively carve out the interpretable models from Black box
2. Concept based, not pixel based
3. First order logic (FOL) for concept interaction

# Assumptions



$f^0$

$\Phi$    $h^0$

$\mathcal{X}$

$\mathcal{C}$

thin ears
shortened muzzle
round feet
....

$\mathcal{Y}$

Dog
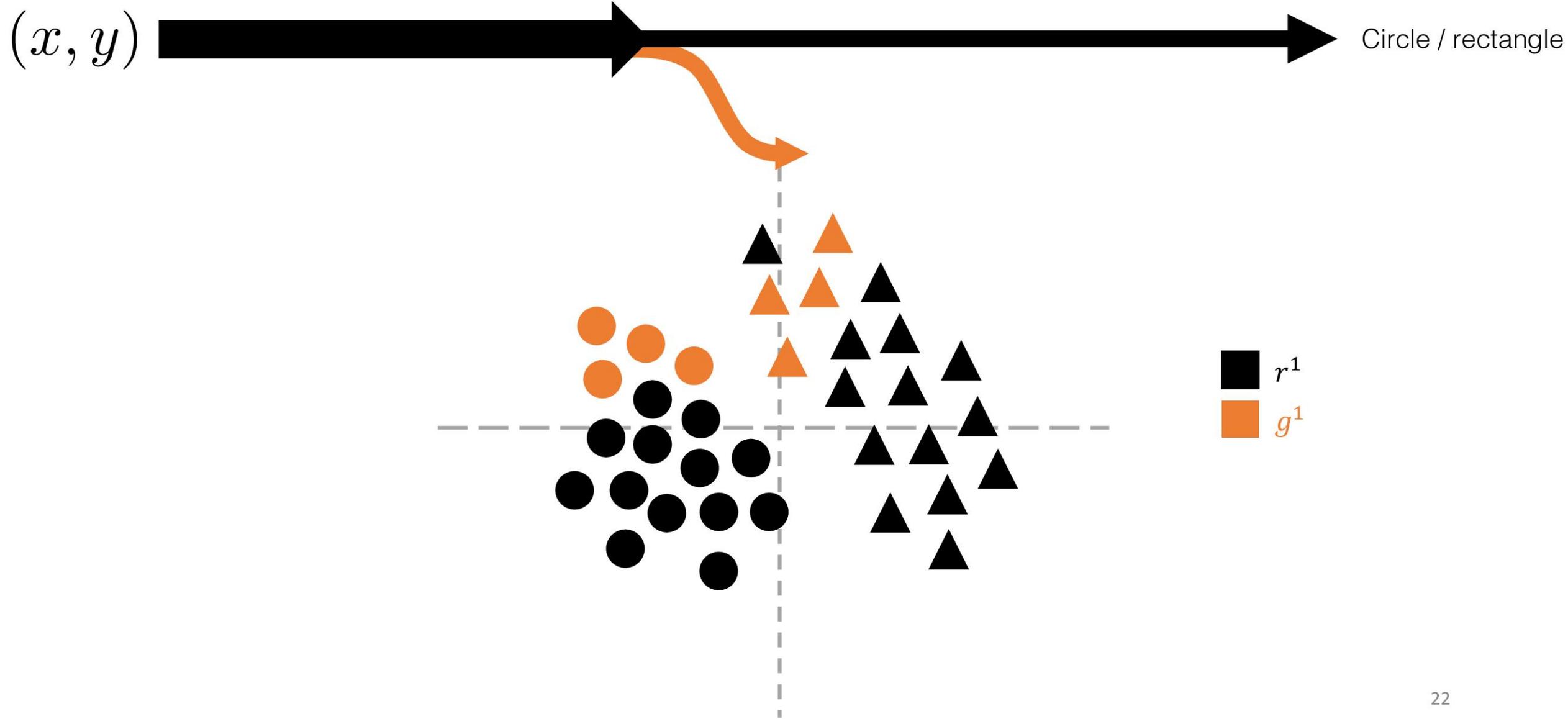
# Iteratively carve out interpretable models

$(x, y)$ ➜ Circle / rectangle



■ $r^0$

# Iteratively carve out interpretable models



$(x, y)$ → → Circle / rectangle

■ $r^1$
■ $g^1$

# Iteratively carve out interpretable models

$(x, y)$

Circle / rectangle

$r^2$

$g^1$

$g^2$

# Iteratively carve out interpretable models

Residual ($r^0$)

Blackbox Model

Interpretable Model

Selector

Each g to produce sample specific FOLs (Barberio et al. AAAI 2022) .

# Iteratively carve out interpretable models



Residual ($r^0$)

Blackbox Model

Interpretable Model

Selector

$\left\{ \; g^1 \; \right\}$

Each g to produce sample specific FOLs (Barberio et al. AAAI 2022) .

# Iteratively carve out interpretable models



Residual ($r^0$)     Residual ($r^1$)

■ Blackbox Model

▬ Interpretable Model $\left\{\begin{array}{c} \square \\ \square\square\square \\ g^1 \end{array}\right\}$

⟳ Selector

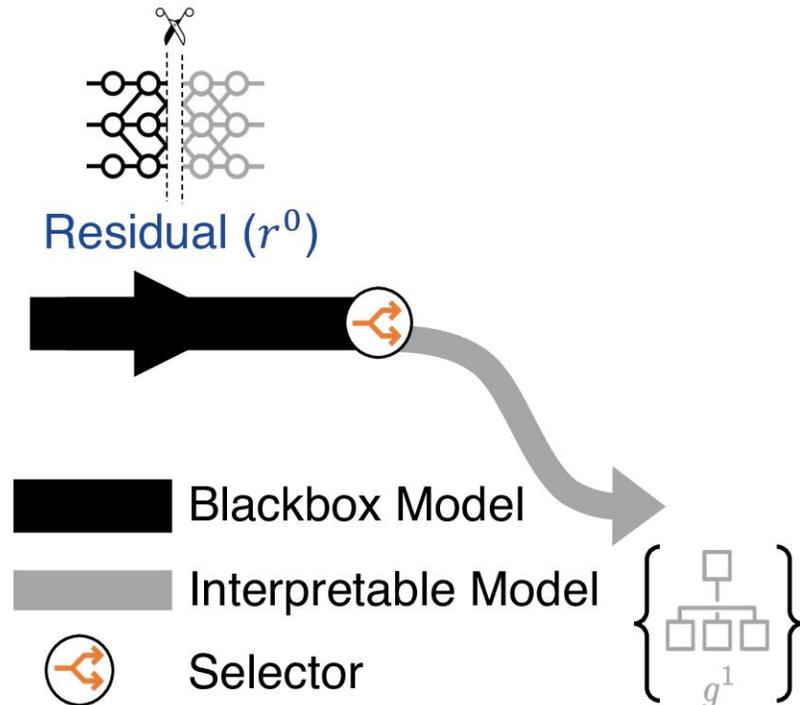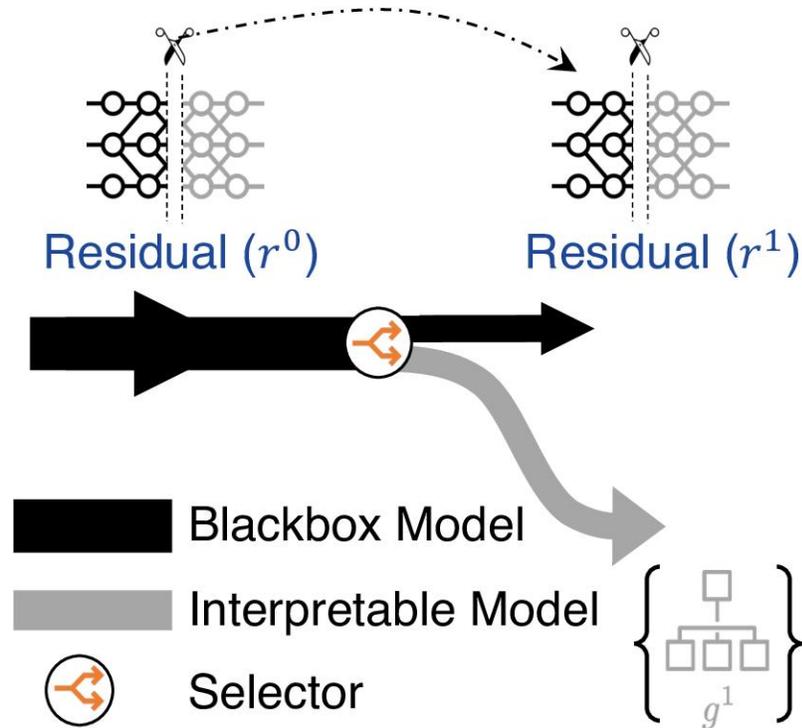Each g to produce sample specific FOLs (Barberio et al. AAAI 2022) .

# Iteratively carve out interpretable models



Each g to produce sample specific FOLs (Barberio et al. AAAI 2022) .
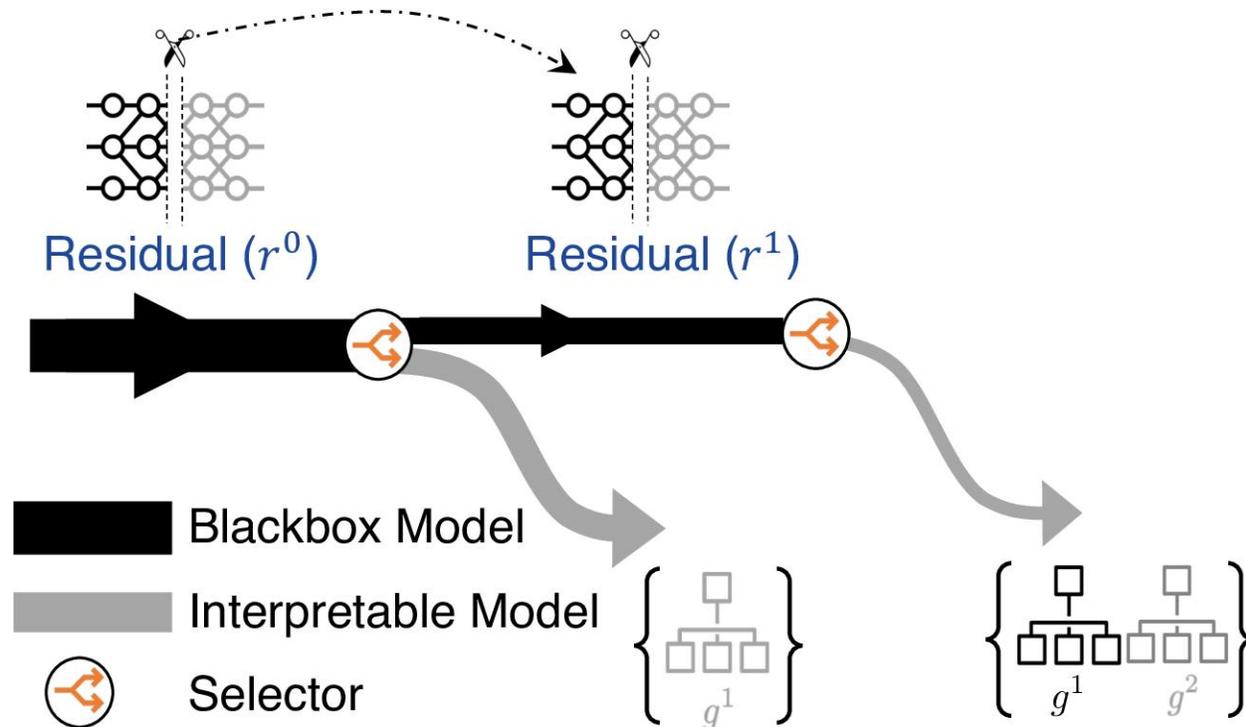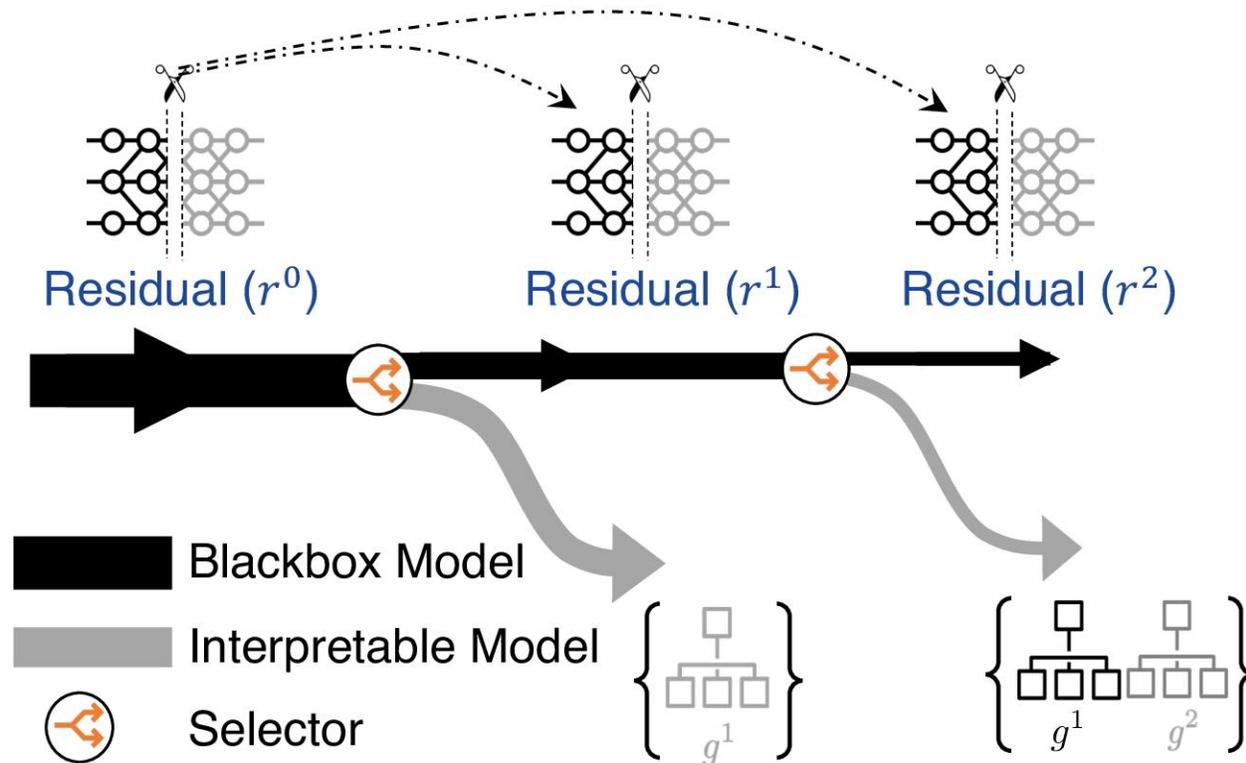
# Iteratively carve out interpretable models



Residual ($r^0$)  Residual ($r^1$)  Residual ($r^2$)

Blackbox Model
Interpretable Model
Selector

$\left\{ \begin{array}{c} \\ g^1 \end{array} \right\}$  $\left\{ \begin{array}{cc} & \\ g^1 & g^2 \end{array} \right\}$

Each g to produce sample specific FOLs (Barberio et al. AAAI 2022) .

# Iteratively carve out interpretable models



Residual ($r^0$)    Residual ($r^1$)    Residual ($r^2$)

Blackbox Model
Interpretable Model
Selector

$\left\{ g^1 \right\}$    $\left\{ g^1 \quad g^2 \right\}$    $\left\{ g^1 \quad g^2 \quad g^3 \right\}$

Each g to produce sample specific FOLs (Barberio et al. AAAI 2022) .

# Iteratively carve out interpretable models



Residual ($r^0$)  Residual ($r^1$)  Residual ($r^2$)  Residual ($r^3$)

Blackbox Model
Interpretable Model
Selector

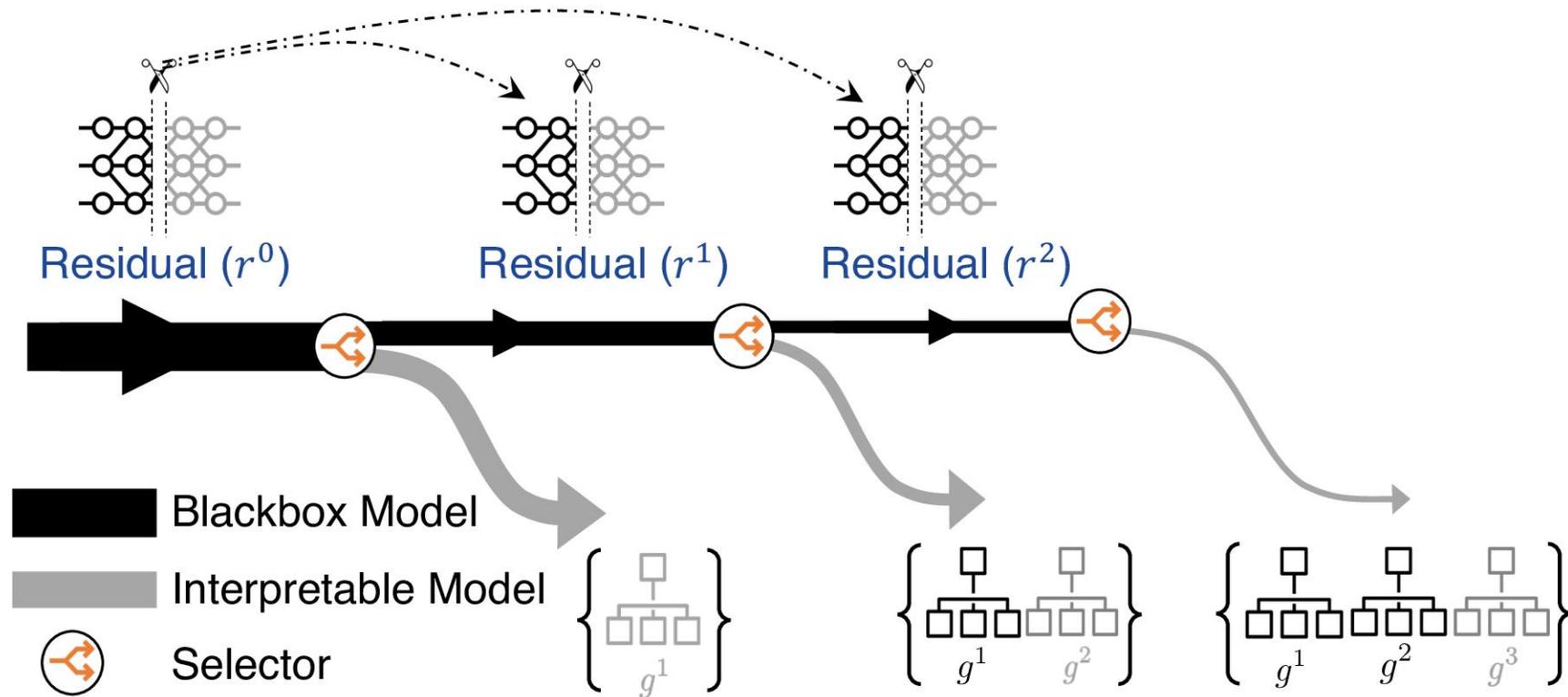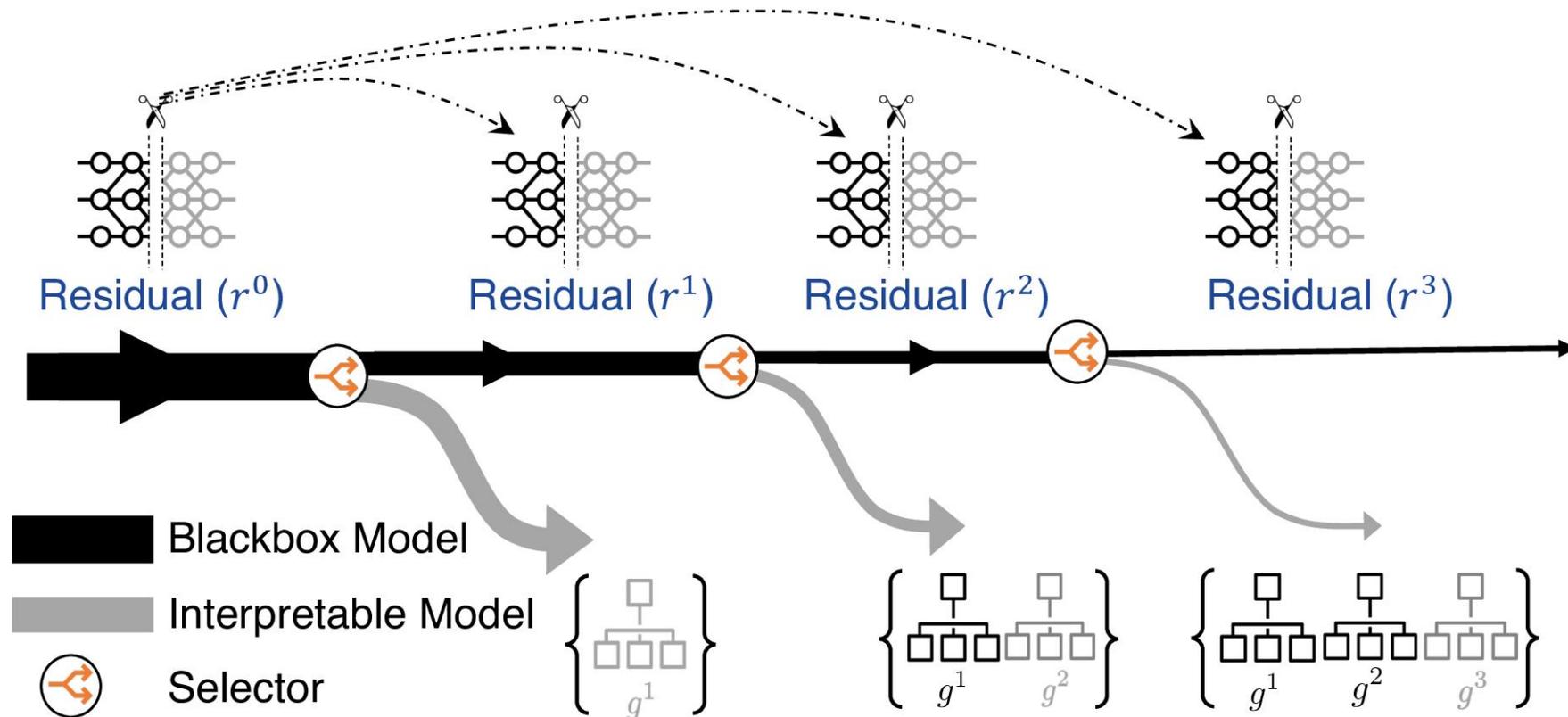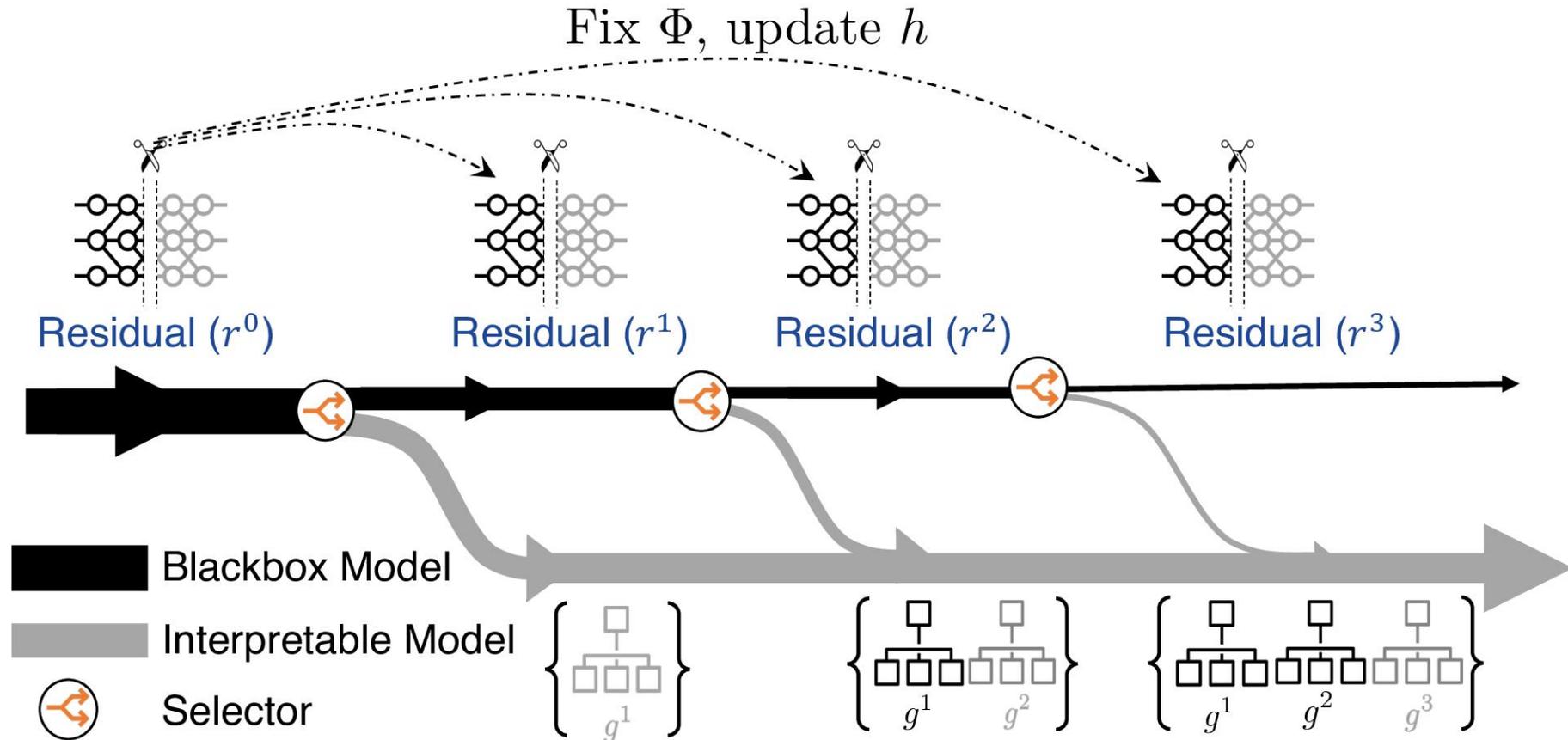$\left\{ g^1 \right\}$  $\left\{ g^1 \quad g^2 \right\}$  $\left\{ g^1 \quad g^2 \quad g^3 \right\}$

Each g to produce sample specific FOLs (Barberio et al. AAAI 2022) .

# Iteratively carve out interpretable models



Each g to produce sample specific FOLs (Barberio et al. AAAI 2022) .