# Initialization Is Critical to Reasoning Ability of Transformer

**Zhongwang Zhang**
**Shanghai Jiao Tong University**
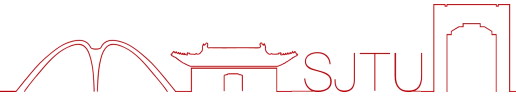**Wednesday, November 13, 2024**

**Try to remember these equations and do test, e.g., (2，1，24) =?**

$(1, 2, 21) = 27, (1, 2, 30) = 36, (1, 2, 47) = 53, (2, 1, 15) = 21, (2, 1, 24) = 30, (2, 1, 41) = 47$

$(3, 4, 10) = 0, (3, 4, 19) = 9, (3, 4, 34) = 24, (3, 4, 41) = 31, (4, 3, 18) = 8, (4, 3, 27) = 17$

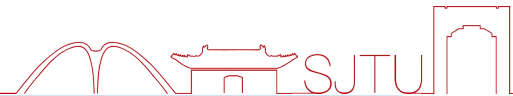$(1, 3, 14) = 17, (1, 3, 18) = 21, (1, 3, 27) = 30, (1, 3, 35) = 38, (3, 1, 9) = 12, (3, 1, 13) = 16$

$(2, 4, 16) = 9, (2, 4, 17) = 10, (2, 4, 30) = 23, (2, 4, 38) = 31, (4, 2, 13) = 6, (4, 2, 14) = 7$

$(1, 4, 32) = 29, (1, 4, 26) = 23, (1, 4, 13) = 10, (1, 4, 6) = 3, (4, 1, 32) = 29, (4, 1, 26) = 23$

$(2, 3, 22) = 21, (2, 3, 17) = 16, (2, 3, 28) = 27, (2, 3, 12) = 11, (3, 2, 22) = 21, (3, 2, 17) = 16$

$\vdots$

# Inference example

$(1, 2, x) = x + 6$
$(2, 1, x) = x + 6$
$(1, 3, x) = x + 3$
$(3, 1, x) = x + 3$
$(1, 4, x) = x - 3$
$(4, 1, x) = x - 3$
$(3, 2, x) = x - 1$
$(2, 2, x) = x - 1$
$(4, 2, x) = x - 9$
$(2, 4, x) = x - 9$
$(3, 4, x) = x + 6$
$(4, 3, x) = x + 6$
$(1, 1, x) = x + 10$
$(2, 2, x) = x + 2$
$(3, 3, x) = x - 4$
$(4, 4, x) = x - 16$

$(1, 2, 21) = 27, \quad (2, 1, 29) = 35$
$(2, 1, 50) = 56, \quad (1, 2, 31) = 37$
$(1, 3, 16) = 19, \quad (3, 1, 16) = 19$
$(1, 4, 67) = 64, \quad (4, 1, 99) = 96$
$(3, 2, 60) = 59, \quad (2, 3, 50) = 49$
$(4, 2, 48) = 39, \quad (2, 4, 33) = 24$
$(3, 4, 77) = 83, \quad (4, 3, 90) = 96$
$(1, 1, 58) = 68, \quad (1, 1, 51) = 61$
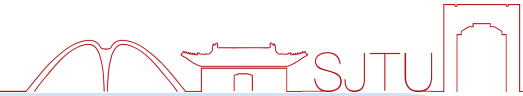$(2, 2, 46) = 48, \quad (2, 2, 35) = 37$
$(3, 3, 36) = 32, \quad (3, 3, 29) = 25$
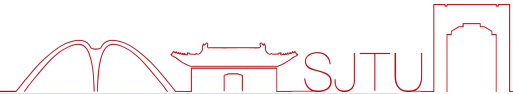$(4, 4, 88) = 54, \quad (4, 4, 46) = 30$

$\vdots$

# Inference example

$$(1, 2, x) = x + 6$$
$$(2, 1, x) = x + 6$$
$$(1, 3, x) = x + 3$$
$$(3, 1, x) = x + 3$$
$$(1, 4, x) = x - 3$$
$$(4, 1, x) = x - 3$$
$$(3, 2, x) = x - 1$$
$$(2, 2, x) = x - 1$$
$$(4, 2, x) = x - 9$$
$$(2, 4, x) = x - 9$$
$$(3, 4, x) = x + 6$$
$$(4, 3, x) = x + 6$$
$$(1, 1, x) = x + 10$$
$$(2, 2, x) = x + 2$$
$$(3, 3, x) = x - 4$$
$$(4, 4, x) = x - 16$$

$$(1, 2, 21) = 27, \quad (2, 1, 29) = 35$$
$$(2, 1, 50) = 56, \quad (1, 2, 31) = 37$$
$$(1, 3, 16) = 19, \quad (3, 1, 16) = 19$$
$$(1, 4, 67) = 64, \quad (4, 1, 99) = 96$$
$$(3, 2, 60) = 59, \quad (2, 3, 50) = 49$$
$$(4, 2, 48) = 39, \quad (2, 4, 33) = 24$$
$$(3, 4, 77) = 83, \quad (4, 3, 90) = 96$$
$$(1, 1, 58) = 68, \quad (1, 1, 51) = 61$$
$$(2, 2, 46) = 48, \quad (2, 2, 35) = 37$$
$$(3, 3, 36) = 32, \quad (3, 3, 29) = 25$$
$$(4, 4, 88) = 54, \quad (4, 4, 46) = 30$$
$$\vdots$$

$$(1, 2, x) = (2, 1, x) = x + 6$$
$$(1, 3, x) = (3, 1, x) = x + 3$$
$$(1, 4, x) = (4, 1, x) = x - 3$$
$$(3, 2, x) = (2, 2, x) = x - 1$$
$$(4, 2, x) = (2, 4, x) = x - 9$$
$$(3, 4, x) = (4, 3, x) = x + 6$$
$$(1, 1, x) = x + 10$$
$$(2, 2, x) = x + 2$$
$$(3, 3, x) = x - 4$$
$$(4, 4, x) = x - 16$$

# Inference example

$$f_1(x) = x + 5$$
$$f_2(x) = x + 1$$
$$f_3(x) = x - 2$$
$$f_4(x) = x - 8$$

$$(1, 2, x) = (2, 1, x) = x + 6$$
$$(1, 3, x) = (3, 1, x) = x + 3$$
$$(1, 4, x) = (4, 1, x) = x - 3$$
$$(3, 2, x) = (2, 2, x) = x - 1$$
$$(4, 2, x) = (2, 4, x) = x - 9$$
$$(3, 4, x) = (4, 3, x) = x + 6$$
$$(1, 1, x) = x + 10$$
$$(2, 2, x) = x + 2$$
$$(3, 3, x) = x - 4$$
$$(4, 4, x) = x - 16$$

$$(1, 2, x) = x + 6$$
$$(2, 1, x) = x + 6$$
$$(1, 3, x) = x + 3$$
$$(3, 1, x) = x + 3$$
$$(1, 4, x) = x - 3$$
$$(4, 1, x) = x - 3$$
$$(3, 2, x) = x - 1$$
$$(2, 2, x) = x - 1$$
$$(4, 2, x) = x - 9$$
$$(2, 4, x) = x - 9$$
$$(3, 4, x) = x + 6$$
$$(4, 3, x) = x + 6$$
$$(1, 1, x) = x + 10$$
$$(2, 2, x) = x + 2$$
$$(3, 3, x) = x - 4$$
$$(4, 4, x) = x - 16$$

$$(1, 2, 21) = 27, \quad (2, 1, 29) = 35$$
$$(2, 1, 50) = 56, \quad (1, 2, 31) = 37$$
$$(1, 3, 16) = 19, \quad (3, 1, 16) = 19$$
$$(1, 4, 67) = 64, \quad (4, 1, 99) = 96$$
$$(3, 2, 60) = 59, \quad (2, 3, 50) = 49$$
$$(4, 2, 48) = 39, \quad (2, 4, 33) = 24$$
$$(3, 4, 77) = 83, \quad (4, 3, 90) = 96$$
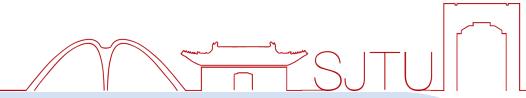$$(1, 1, 58) = 68, \quad (1, 1, 51) = 61$$
$$(2, 2, 46) = 48, \quad (2, 2, 35) = 37$$
$$(3, 3, 36) = 32, \quad (3, 3, 29) = 25$$
$$(4, 4, 88) = 54, \quad (4, 4, 46) = 30$$
$$\vdots$$

# Inference example

**Inference Complexity**

**Memory Complexity**

$$f_1(x) = x + 5$$
$$f_2(x) = x + 1$$
$$f_3(x) = x - 2$$
$$f_4(x) = x - 8$$

$$(1, 2, x) = (2, 1, x) = x + 6$$
$$(1, 3, x) = (3, 1, x) = x + 3$$
$$(1, 4, x) = (4, 1, x) = x - 3$$
$$(3, 2, x) = (2, 2, x) = x - 1$$
$$(4, 2, x) = (2, 4, x) = x - 9$$
$$(3, 4, x) = (4, 3, x) = x + 6$$
$$(1, 1, x) = x + 10$$
$$(2, 2, x) = x + 2$$
$$(3, 3, x) = x - 4$$
$$(4, 4, x) = x - 16$$

$$(1, 2, x) = x + 6$$
$$(2, 1, x) = x + 6$$
$$(1, 3, x) = x + 3$$
$$(3, 1, x) = x + 3$$
$$(1, 4, x) = x - 3$$
$$(4, 1, x) = x - 3$$
$$(3, 2, x) = x - 1$$
$$(2, 2, x) = x - 1$$
$$(4, 2, x) = x - 9$$
$$(2, 4, x) = x - 9$$
$$(3, 4, x) = x + 6$$
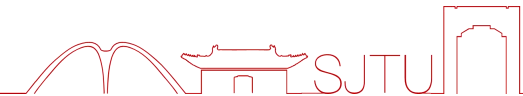$$(4, 3, x) = x + 6$$
$$(1, 1, x) = x + 10$$
$$(2, 2, x) = x + 2$$
$$(3, 3, x) = x - 4$$
$$(4, 4, x) = x - 16$$

$$(1, 2, 21) = 27, \quad (2, 1, 29) = 35$$
$$(2, 1, 50) = 56, \quad (1, 2, 31) = 37$$
$$(1, 3, 16) = 19, \quad (3, 1, 16) = 19$$
$$(1, 4, 67) = 64, \quad (4, 1, 99) = 96$$
$$(3, 2, 60) = 59, \quad (2, 3, 50) = 49$$
$$(4, 2, 48) = 39, \quad (2, 4, 33) = 24$$
$$(3, 4, 77) = 83, \quad (4, 3, 90) = 96$$
$$(1, 1, 58) = 68, \quad (1, 1, 51) = 61$$
$$(2, 2, 46) = 48, \quad (2, 2, 35) = 37$$
$$(3, 3, 36) = 32, \quad (3, 3, 29) = 25$$
$$(4, 4, 88) = 54, \quad (4, 4, 46) = 30$$
$$\vdots$$

# Reasoning example

**Reasoning Complexity**

**Memory Complexity**

**Memorizing**

## Composition

$$f_1(x) = x + 5$$
$$f_2(x) = x + 1$$
$$f_3(x) = x - 2$$
$$f_4(x) = x - 8$$

## Symmetry

$$(1, 2, x) = (2, 1, x) = x + 6$$
$$(1, 3, x) = (3, 1, x) = x + 3$$
$$(1, 4, x) = (4, 1, x) = x - 3$$
$$(3, 2, x) = (2, 2, x) = x - 1$$
$$(4, 2, x) = (2, 4, x) = x - 9$$
$$(3, 4, x) = (4, 3, x) = x + 6$$
$$(1, 1, x) = x + 10$$
$$(2, 2, x) = x + 2$$
$$(3, 3, x) = x - 4$$
$$(4, 4, x) = x - 16$$

## Function relation

$$(1, 2, x) = x + 6$$
$$(1, 3, x) = x + 3$$
$$(3, 1, x) = x + 3$$
$$(1, 4, x) = x - 3$$
$$(4, 1, x) = x - 3$$
$$(3, 2, x) = x - 1$$
$$(2, 2, x) = x - 1$$
$$(4, 2, x) = x - 9$$
$$(2, 4, x) = x - 9$$
$$(3, 4, x) = x + 6$$
$$(4, 3, x) = x + 6$$
$$(1, 1, x) = x + 10$$
$$(2, 2, x) = x + 2$$
$$(3, 3, x) = x - 4$$
$$(4, 4, x) = x - 16$$

## Memorizing

$$(1, 2, 21) = 27, \quad (2, 1, 29) = 35$$
$$(2, 1, 50) = 56, \quad (1, 2, 31) = 37$$
$$(1, 3, 16) = 19, \quad (3, 1, 16) = 19$$
$$(1, 4, 67) = 64, \quad (4, 1, 99) = 96$$
$$(3, 2, 60) = 59, \quad (2, 3, 50) = 49$$
$$(4, 2, 48) = 39, \quad (2, 4, 33) = 24$$
$$(3, 4, 77) = 83, \quad (4, 3, 90) = 96$$
$$(1, 1, 58) = 68, \quad (1, 1, 51) = 61$$
$$(2, 2, 46) = 48, \quad (2, 2, 35) = 37$$
$$(3, 3, 36) = 32, \quad (3, 3, 29) = 25$$
$$(4, 4, 88) = 54, \quad (4, 4, 46) = 30$$
$$\vdots$$

上海交通大學
SHANGHAI JIAO TONG UNIVERSITY

# Anchor function

Designated Token
as **Anchor**

| | | |
|---|---|---|
| **1** | : | **+5** |
| **2** | : | **+1** |
| **3** | : | **-2** |
| **4** | : | **-8** |

Other Token
as **Key**

**5**

**6**

⋮

**100**

Data set

(1, 2, 21, 27)

(2, 1, 29, 35)

(2, 1, 50, 56)

⋮

(3, 3, 29, 25)

(4, 4, 46, 30)

**Anchor    Key    Target**

**Train model to
predict results**

Data set

(**1**, **2**, **21**, 27)

(**2**, **1**, **29**, 35)

(**2**, **1**, **50**, 56)

⋮

(**3**, **3**, **29**, 25)

(**4**, **4**, **46**, 30)

**Anchor**   **Key**   **Target**

Loss only here

Calculate loss

Output   Target

NN

Designated Token
as **Anchor**

Other Token
as **Key**

# Composite Anchor function

# Can it learn [4,3]?

# Yes, it can learn [4,3]!

# Symmetric or Inferential?

# Composite Anchor function



**14 seen inferential composite anchors**

| 1 | 1 | : | +10 |

| 1 | 2 | : | +6 |

. . . .

**Composition** →

| 4 | 4 | : | -16 |

**Padding** →

**1 seen non-inferential composite anchors**

| 3 | 4 | : | -6 |

**1 unseen composite anchor**

| 4 | 3 | : | ? |

Left column:
| 1 | : | +5 |
| 2 | : | +1 |
| 3 | : | -2 |
| 4 | : | -8 |

# Phase diagram of symmetric solution



Mechanism 1: learn symmetric structure

$$\boxed{a\ b} = \boxed{b\ a}$$
$$\downarrow$$
$$\boxed{4\ 3} = \boxed{3\ 4} \quad : \quad \boxed{-6}$$

$$\text{Initialization} \sim N\left(0, \frac{1}{d_{in}^{\gamma}}\right)$$

(a)

Large ini    Small ini

acc of symmetric solution with anchor 43

initialization rate γ

bad generalization on seen anchors (test accuracy < 90%)

# Phase diagram of inferential solution



Initialization$\sim N\left(0, \dfrac{1}{d_{in}^{\gamma}}\right)$

Mechanism 2: infer single anchor mappings

$3$ : $-2$  $4$ : $-8$

composite form the infered single anchors

$4\ 3$ : $-10$

(b)

Large ini

Small ini

acc of inferential solution with anchor 43

initialization rate $\gamma$

layer number

bad generalization on seen anchors (test accuracy < 90%)

# Condensation of $W^{Q(1)}$ by column

$$Q^{(l)} = X^{(l)} W^{Q(l)}$$

Cosine Similarity b/w

Columns of $W^{Q(1)}$

Input weight of Q neurons

$X^{(1)}$

$W^{Q(1)}$    output: $Q^{(1)}$

Small ini: clear condensation      Large ini: no condensation

**T-SNE**

Small init: clear structure

Large init: no structure
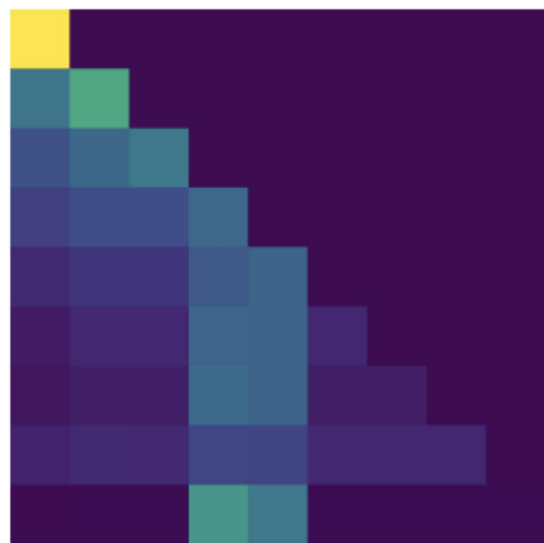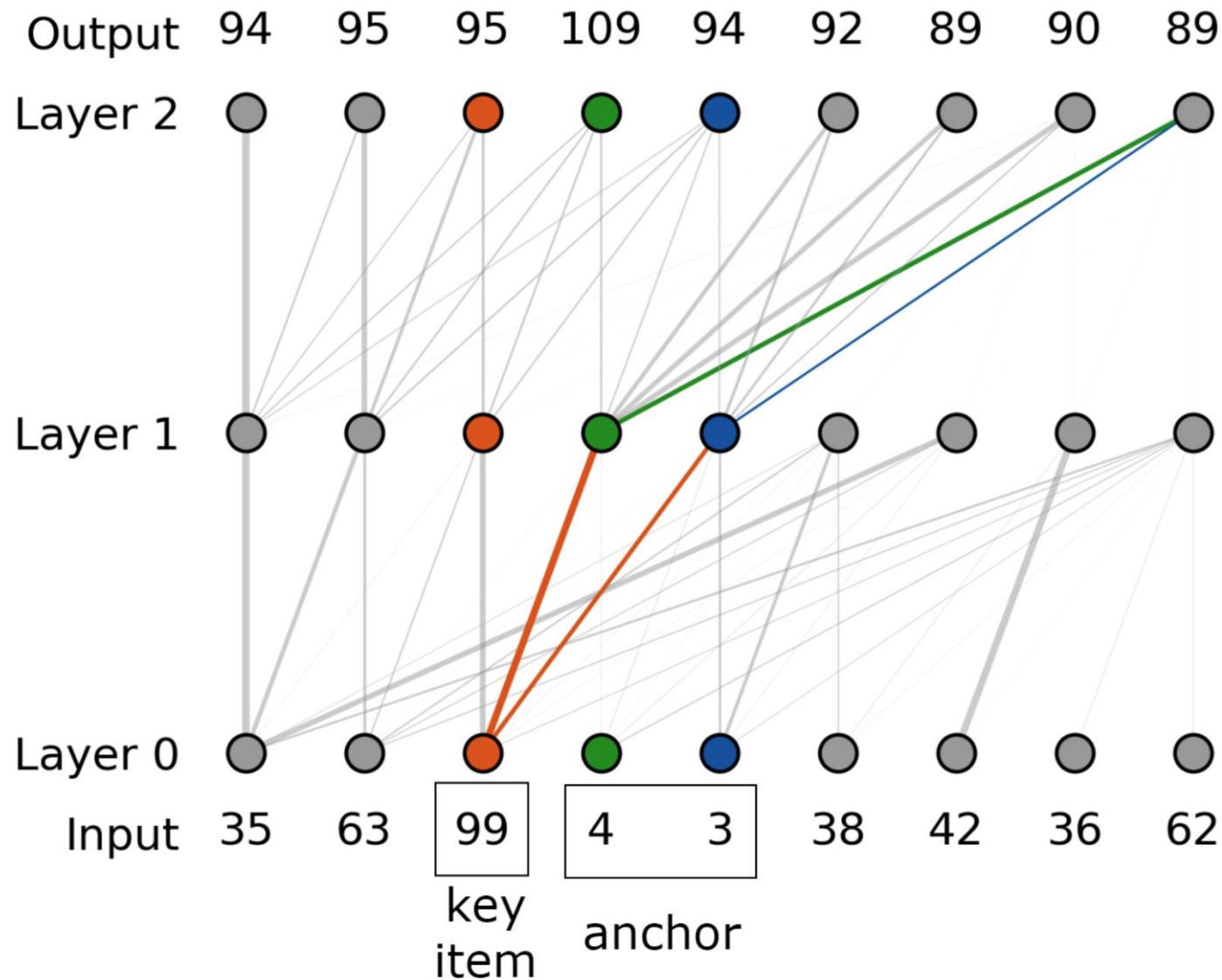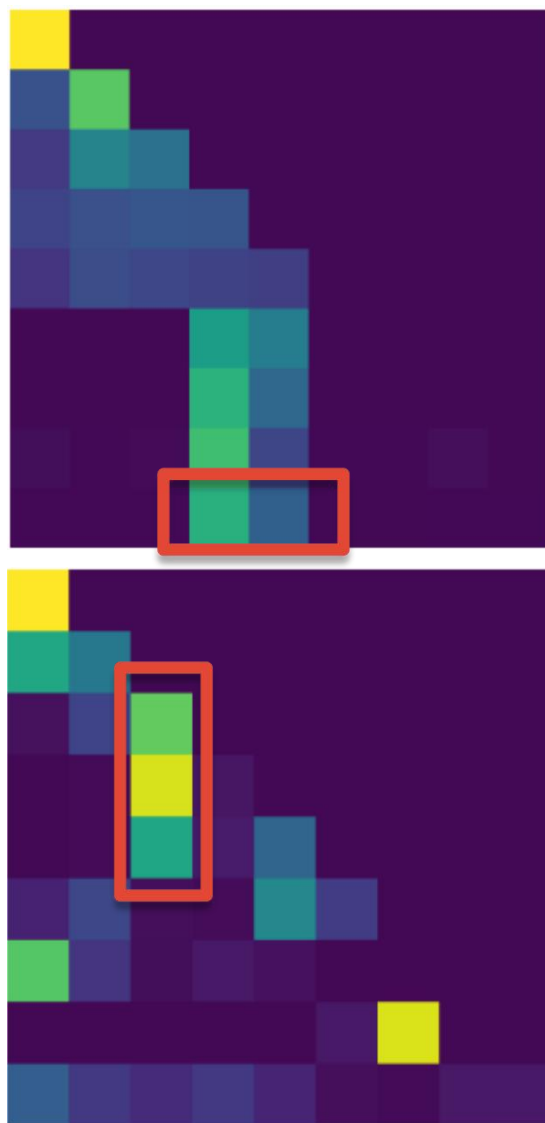
$$f_a(x) = \begin{cases} x + 5, & if\ a\ = 1 \\ x + 1, & if\ a\ = 2 \\ x - 2, & if\ a\ = 3 \\ x - 8, & if\ a\ = 4 \end{cases}$$
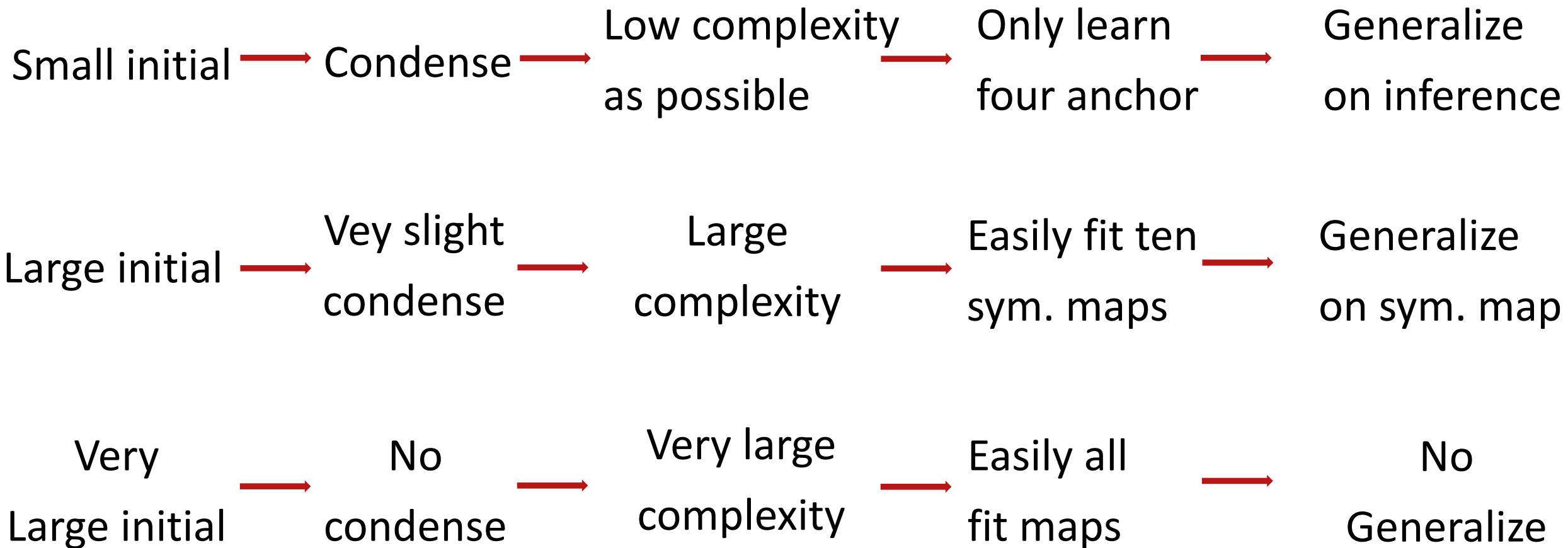
# Symmetric solution

# Inferential solution

# Mechanisms underlying initialization effect

Small initial $\longrightarrow$ Condense $\longrightarrow$ Low complexity as possible $\longrightarrow$ Only learn four anchor $\longrightarrow$ Generalize on inference

Large initial $\longrightarrow$ Vey slight condense $\longrightarrow$ Large complexity $\longrightarrow$ Easily fit ten sym. maps $\longrightarrow$ Generalize on sym. map

Very Large initial $\longrightarrow$ No condense $\longrightarrow$ Very large complexity $\longrightarrow$ Easily all fit maps $\longrightarrow$ No Generalize

上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

# Other realistic reasoning dataset

Large initial ⟶ Large complexity ⟶ Memorization to inference

Small initial ⟶ Low complexity as possible ⟶ Inference

More data, more steps to memorize

上海交通大学

SHANGHAI JIAO TONG UNIVERSITY