



Practical Serializability Verification

Kyle Kingsbury

&

Peter Alvaro

HIGH PERFORMANCE

TRANSACTION

CYBERSECURITY

TRANSACTION

TRANSACTION

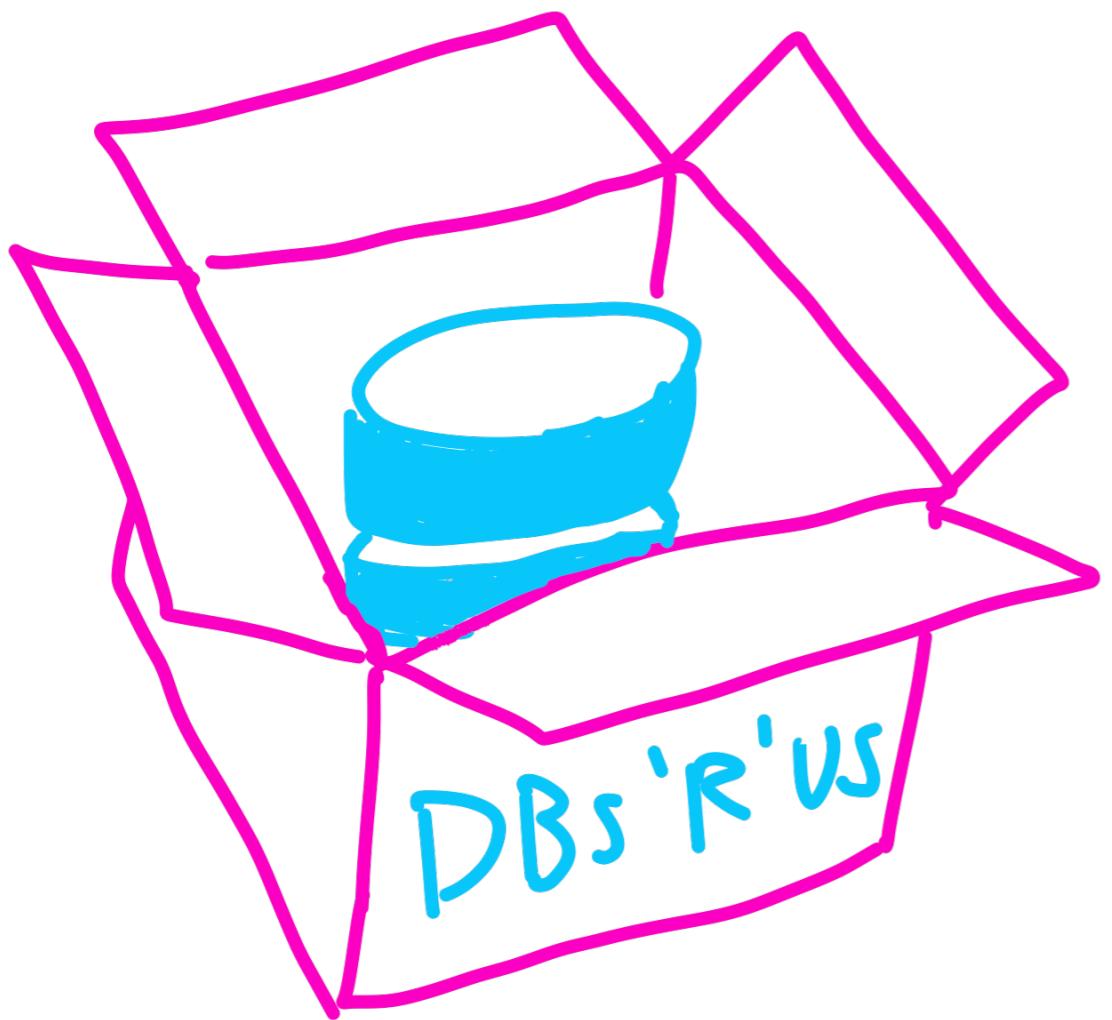
TRANSACTION

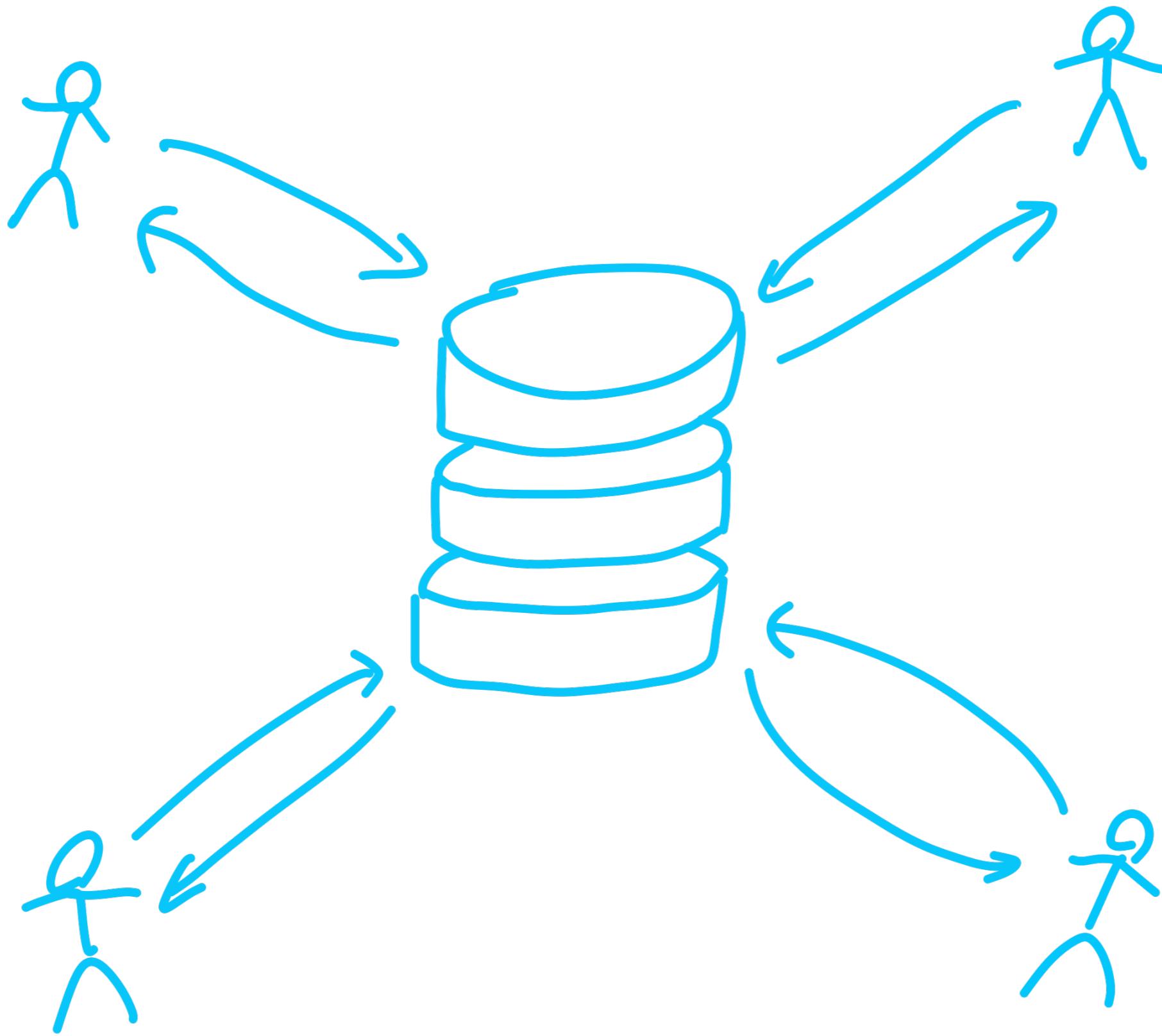
TRANSACTION

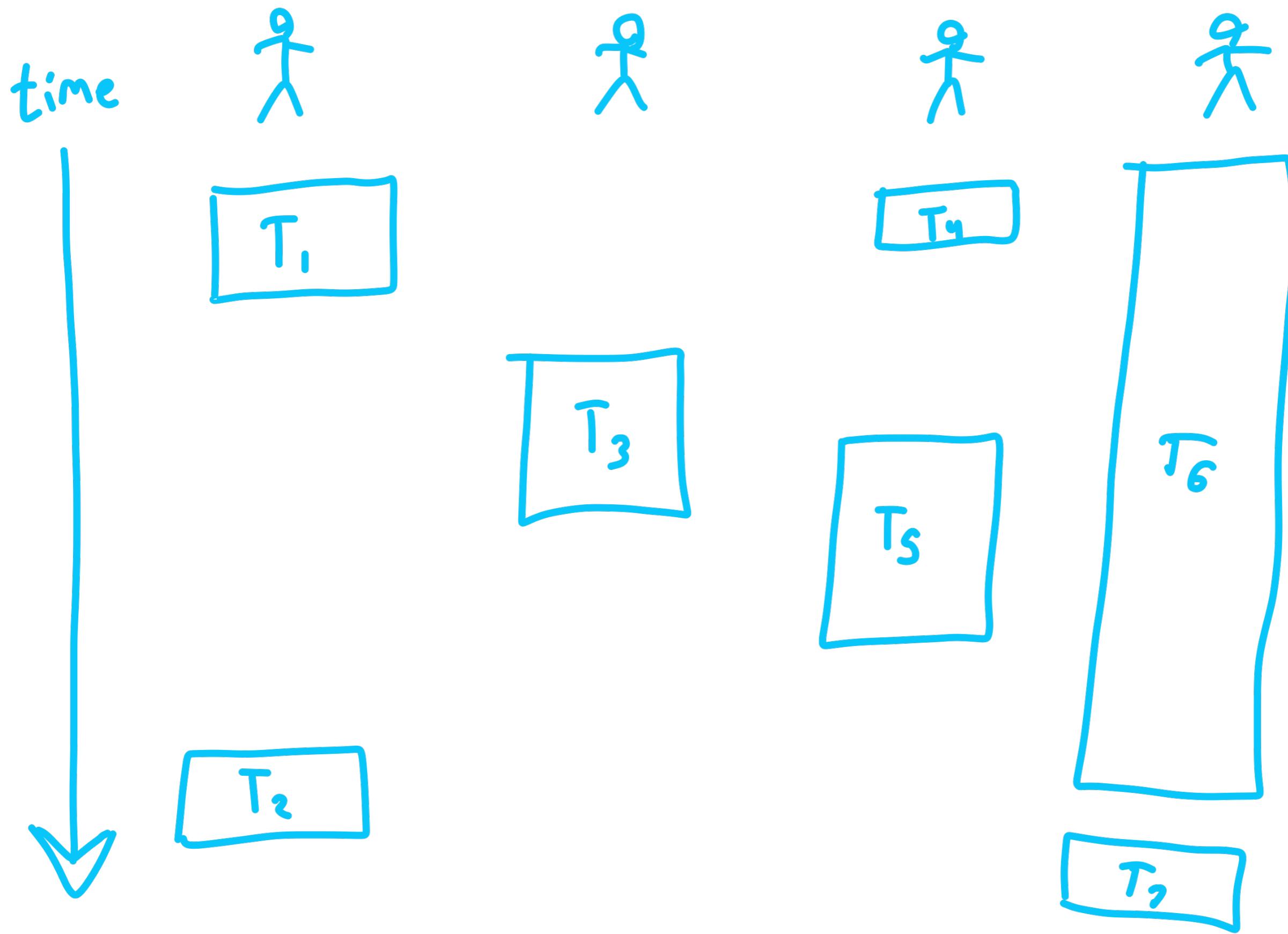
TRANSACTION

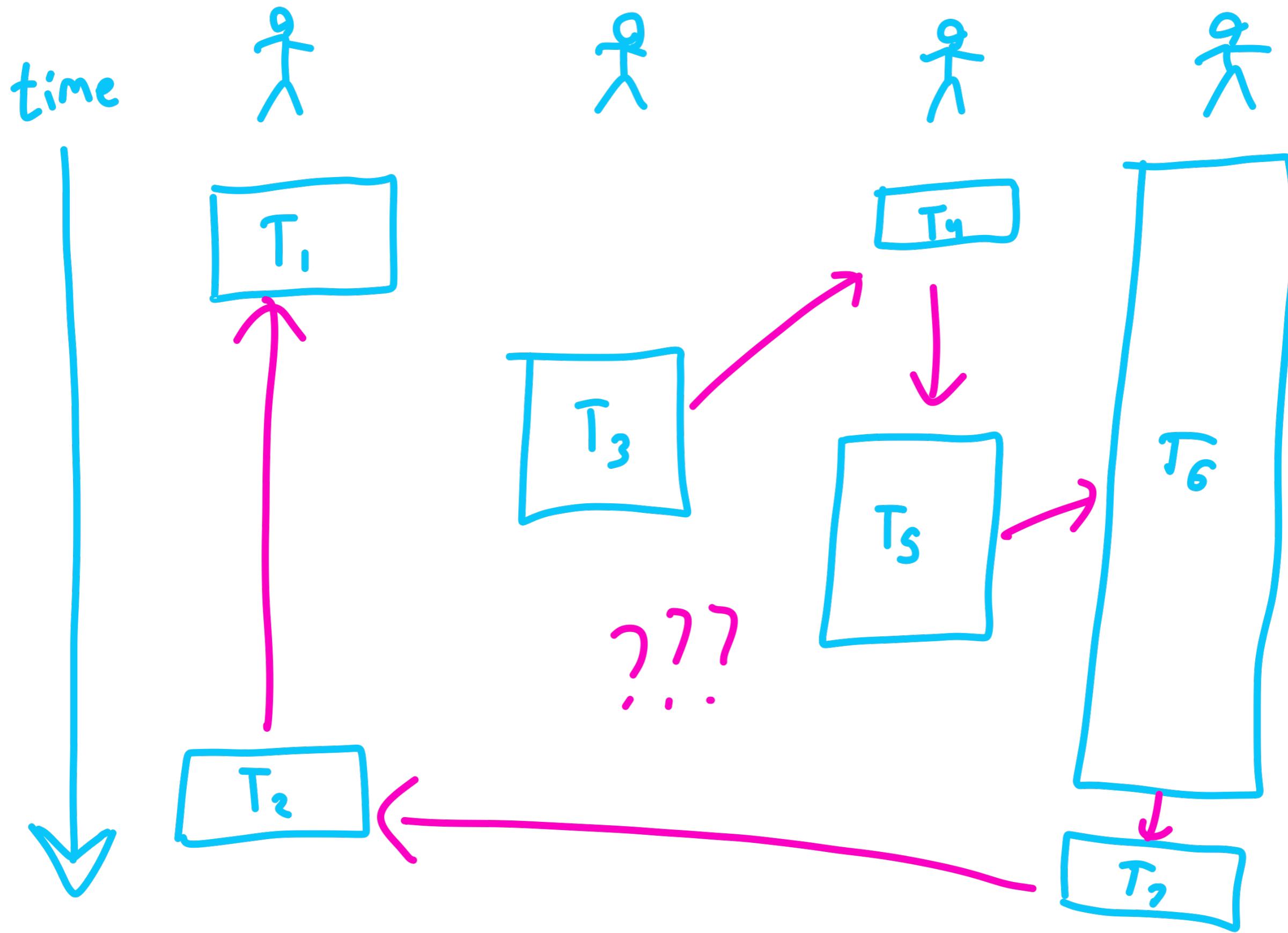
TRANSACTION

So, is it
Serializable?









~~T₁ T₂ T₃ T₄ T₅ T₆ T₇~~

~~T₂ T₁ T₃ T₄ T₅ T₆ T₇~~

~~T₁ T₃ T₂ T₄ T₅ T₆ T₇~~

~~T₂ T₃ T₁ T₄ T₅ T₆ T₇~~

~~T₃ T₁ T₂ T₄ T₅ T₆ T₇~~

~~T₃ T₂ T₁ T₄ T₅ T₆ T₇~~

~~T₁ T₂ T₄ T₃ T₅ T₆ T₇~~

~~T₁ T₄ T₂ T₃ T₅ T₆ T₇~~

~~T₁ T₄ T₃ T₂ T₅ T₆ T₇~~

Aha! Serializable!

$T_1 T_2 T_3 T_4 \bar{T}_5 \bar{T}_6 T_7$
 $\bar{T}_2 T_1 T_3 T_4 \bar{T}_5 \bar{T}_6 T_7$
 $T_1 T_3 \bar{T}_2 T_4 \bar{T}_5 \bar{T}_6 T_7$
 $\bar{T}_2 T_3 T_1 T_4 \bar{T}_5 \bar{T}_6 T_7$
 $T_3 T_1 \bar{T}_2 T_4 \bar{T}_5 \bar{T}_6 T_7$
 $T_3 T_2 T_1 T_4 \bar{T}_5 \bar{T}_6 T_7$
 $T_1 T_2 T_4 T_3 \bar{T}_5 \bar{T}_6 T_7$
 $T_1 T_4 T_2 T_3 \bar{T}_5 \bar{T}_6 T_7$
 $T_1 T_4 T_3 \bar{T}_2 \bar{T}_5 \bar{T}_6 T_7$

N!

Gretchen

github.com/aphyr/gretchen

(Cerone, Bernardi, Gotsman 2015)

Serializable



Check in
each txn
separately

Build read
dependency
constraint
problem in
CNF

Solve for order
using integer
constraint solver

$T_1 : w(x, 3)$  $T_2 : r(x, 3)$ $T_1 < T_2$

$T_1 : w(x, 3)$ $T_2 : r(x, 3)$ $\bar{T}_3 : w(x, 3)$ 

$$(T_1 < T_2) \vee (\bar{T}_3 < \bar{T}_2)$$

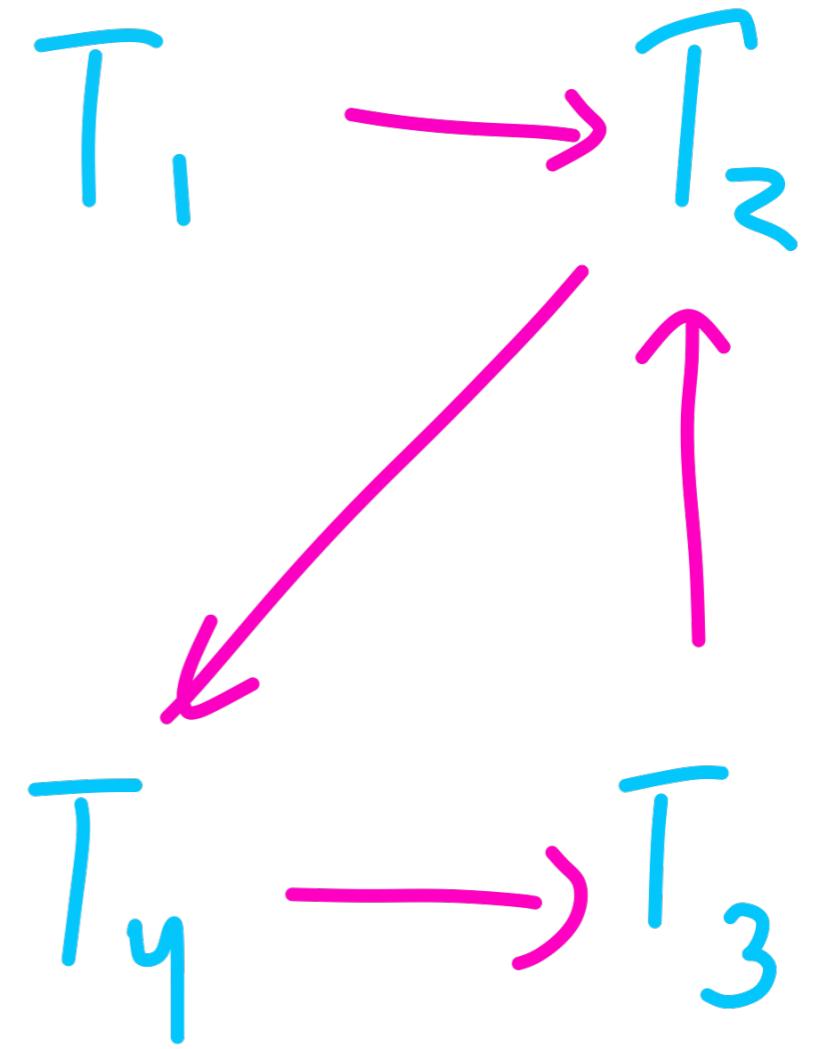
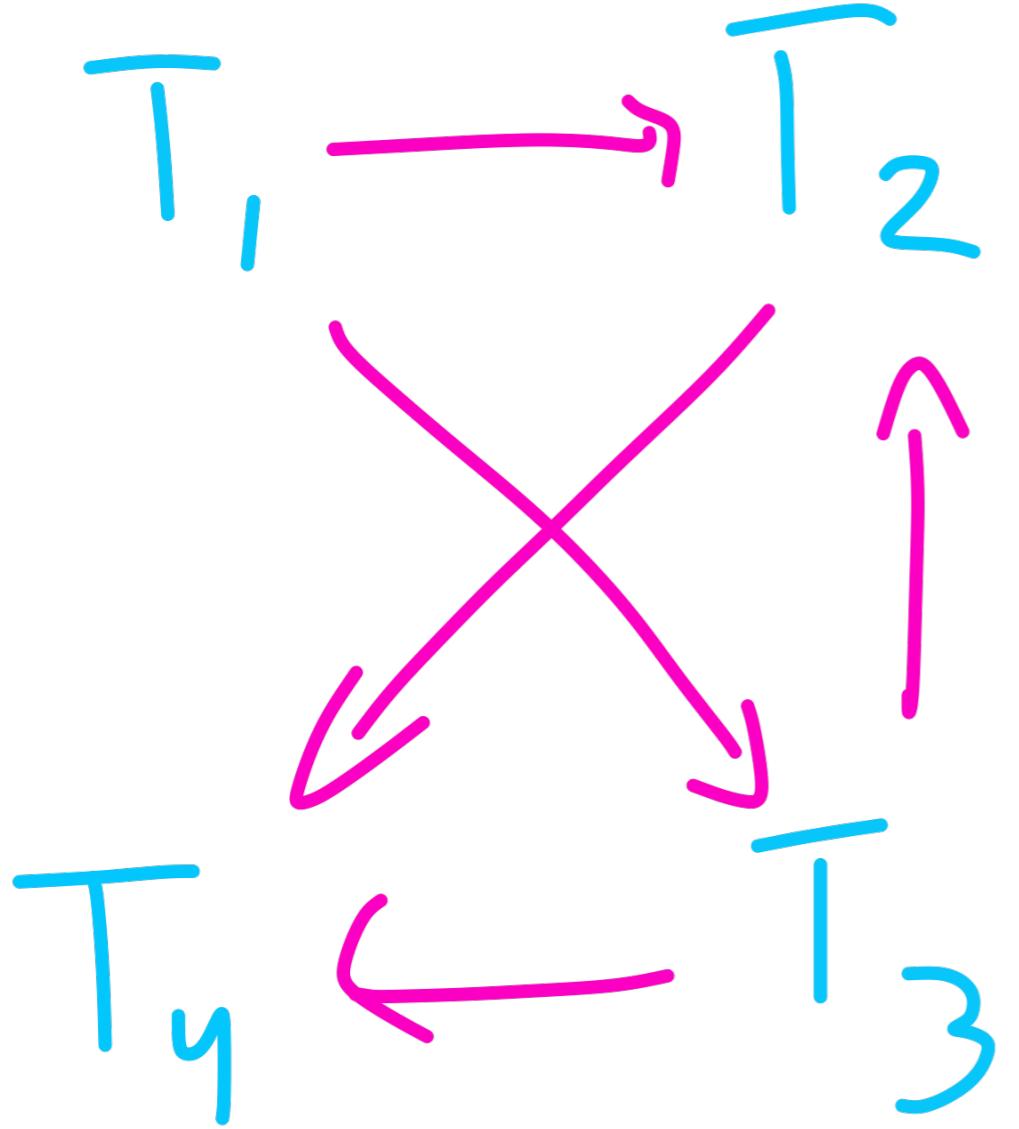
- SAT is ALSO NP-Complete
- Breaks down after ~ 100 txns
- Doesn't explain why no soln!

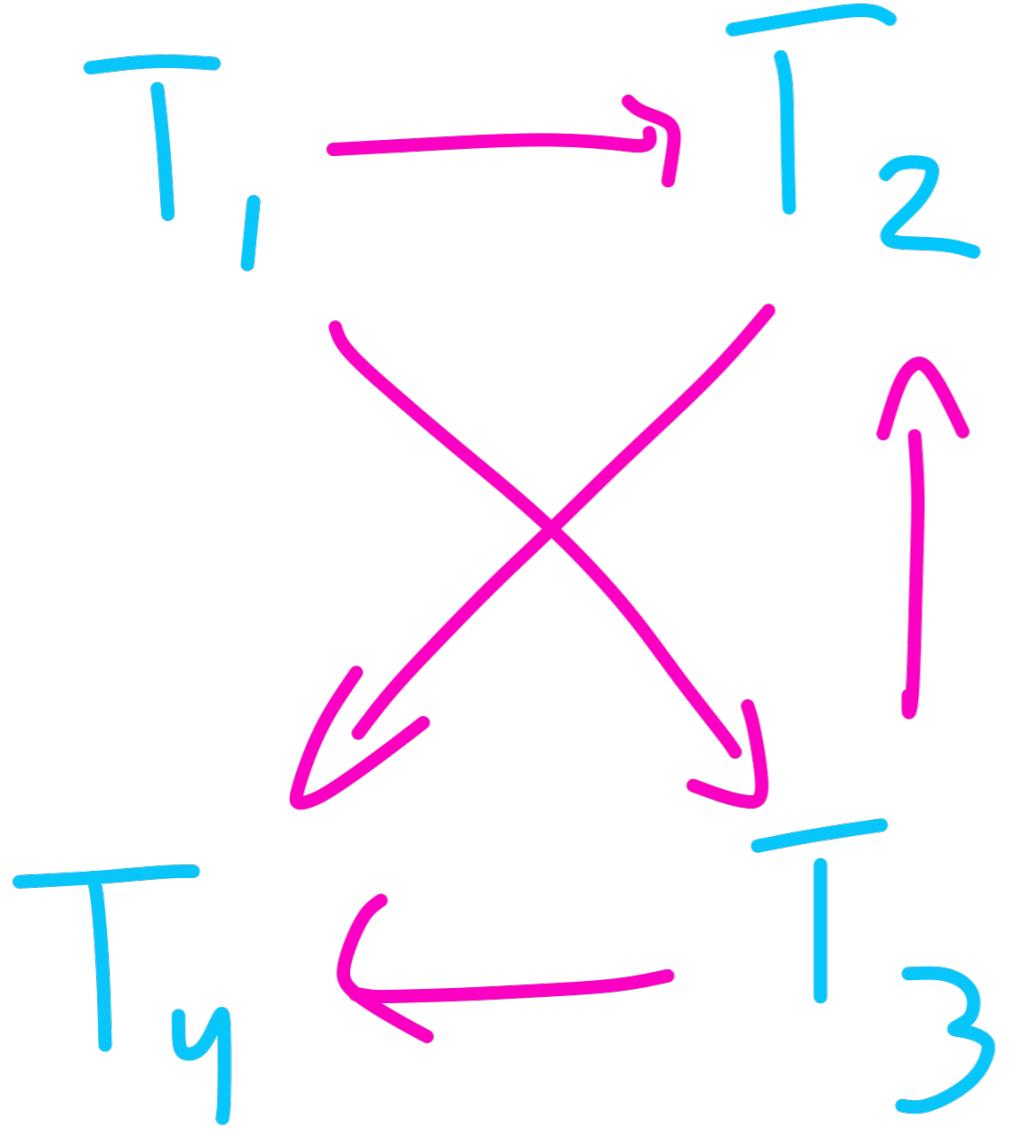


Cycle Detection

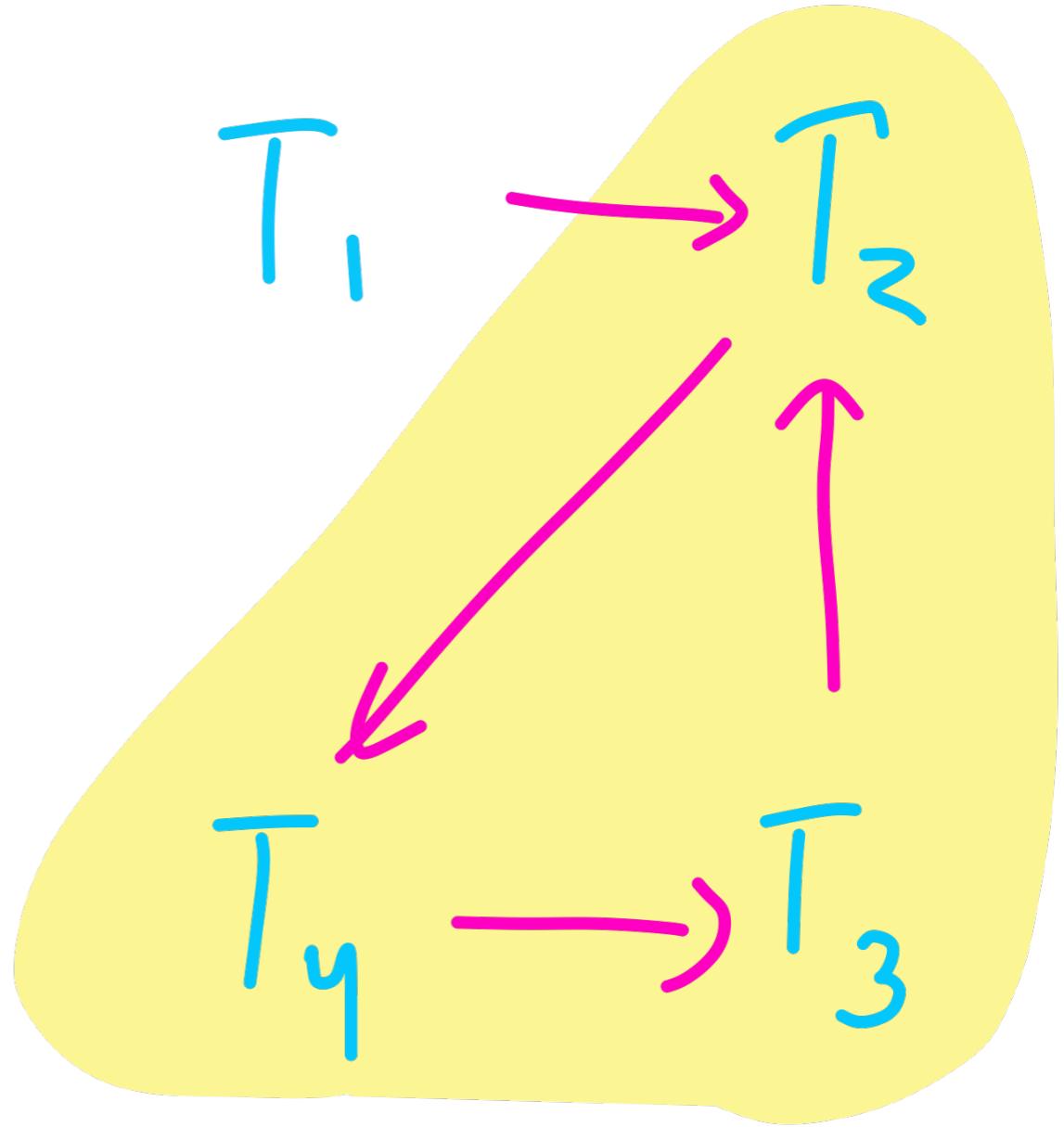


Possibility of a Total Order





Serializable



Not Serializable!

Adya 1995

Isolation Levels in terms
of Dependency Cycles

Write - Read Dependency

$T_1: w(x_i) \rightarrow T_2: r(x_i)$

Write - Write Dependency

$T_1: w(x_i) \rightarrow T_2: w(x_j)$

$x_i \ll \underline{x_j}$

Version Order \ll

Total order of installed xs, ys, etc.

Read - Write Dependency

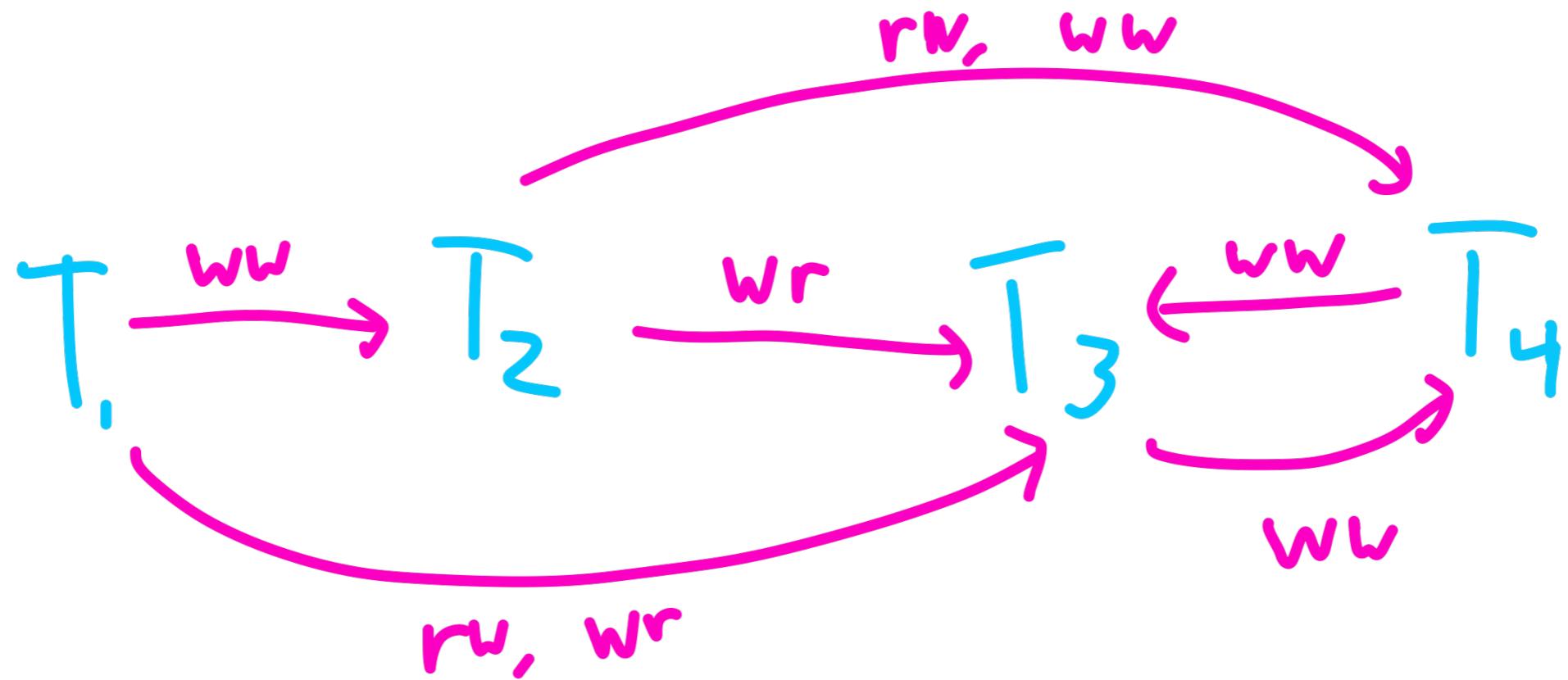
$T_1: r(x_i) \rightarrow T_2: w(x_j)$

$x_i \ll \underline{x_j}$

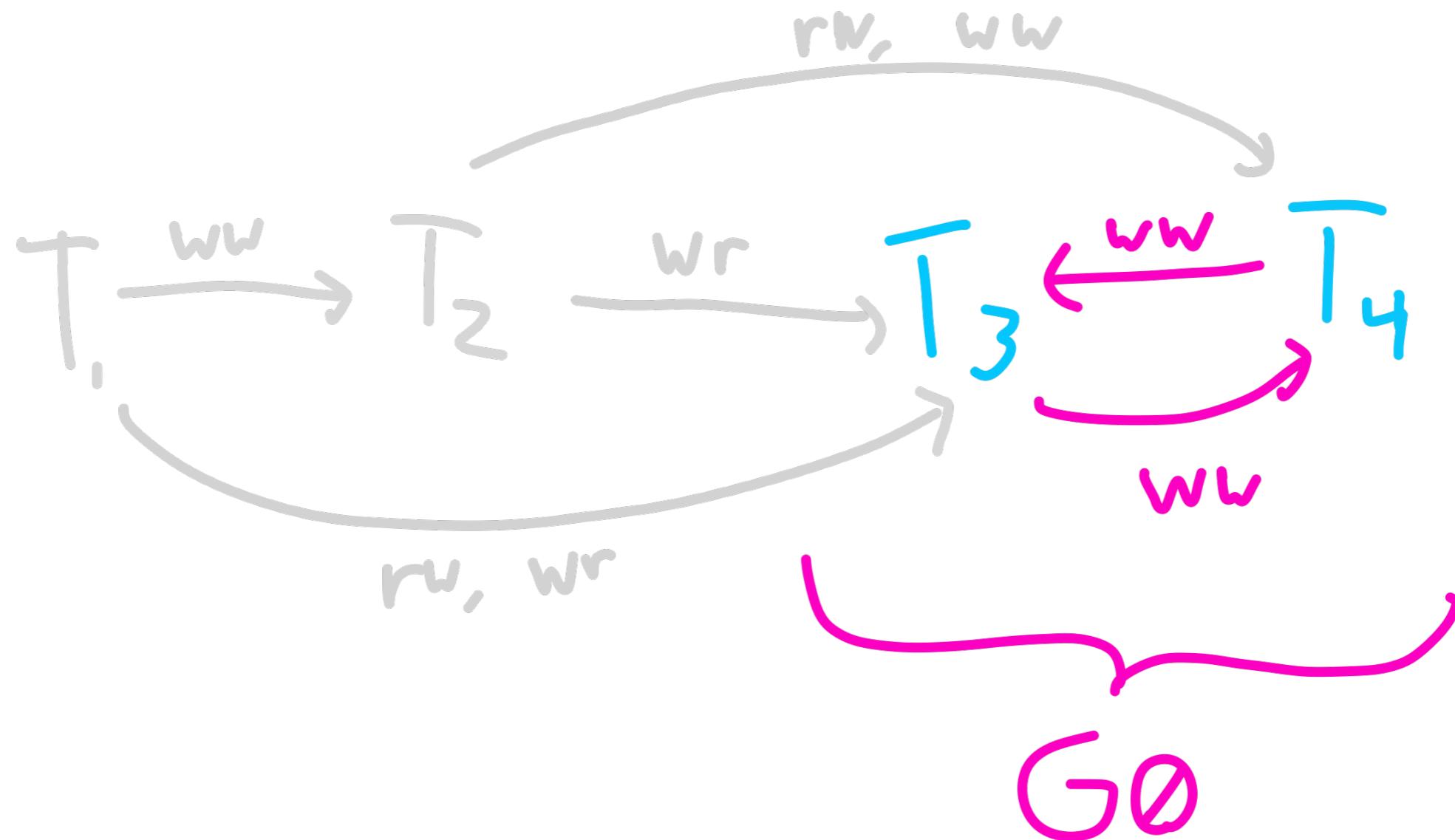
Version Order \ll

Total order of installed xs, ys, etc.

Direct Serialization Graph



Direct Serialization Graph



G0

ww cycle

G1a

aborted read

G1b

intermediate read

G1c

ww + wr cycle

G-single

ww+wr+1rw cycle

G2

ww+wr+rw cycle

Great! Where do we get
an Adya History?

$\bar{T}_1: w(x_1), w(y_1)$

$T_2: r(x_1), w(y_2)$

Great! Where do we get
an Adya History?

$T_1: w(x_1), w(y_1)$ Missing
 ↓? ↓? ↑?
 $y_1 \ll y_2$
 $T_2: r(x_1), w(y_2)$ or
 Only write of x_1 ?
 $y_2 \ll y_1$

$T_1 : w(x, 3)$ $T_2 : r(x, 3)$ ~~$T_3 : w(x, 3)$~~

Recoverability

Blind Writes Destroy History

$w(x_i)$

$w(x_j)$

$w(x_k)$

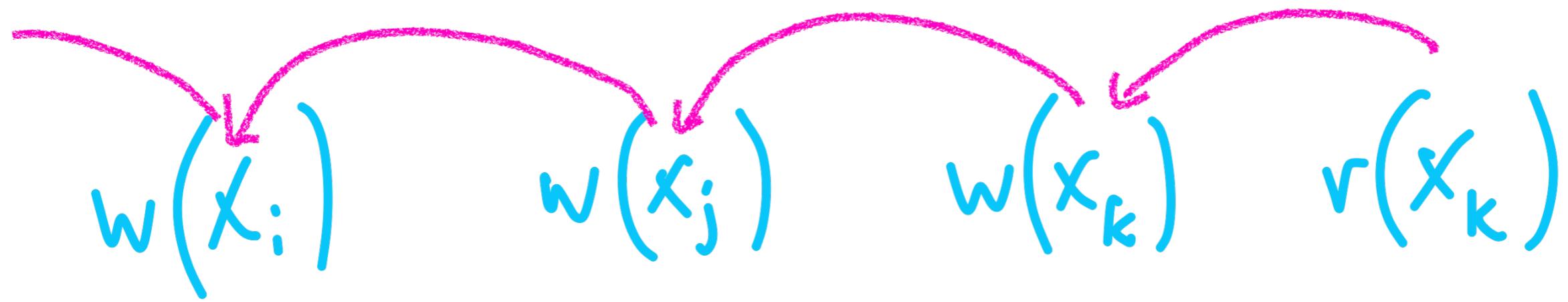
$r(x_k)$

?

?

?





Traceability

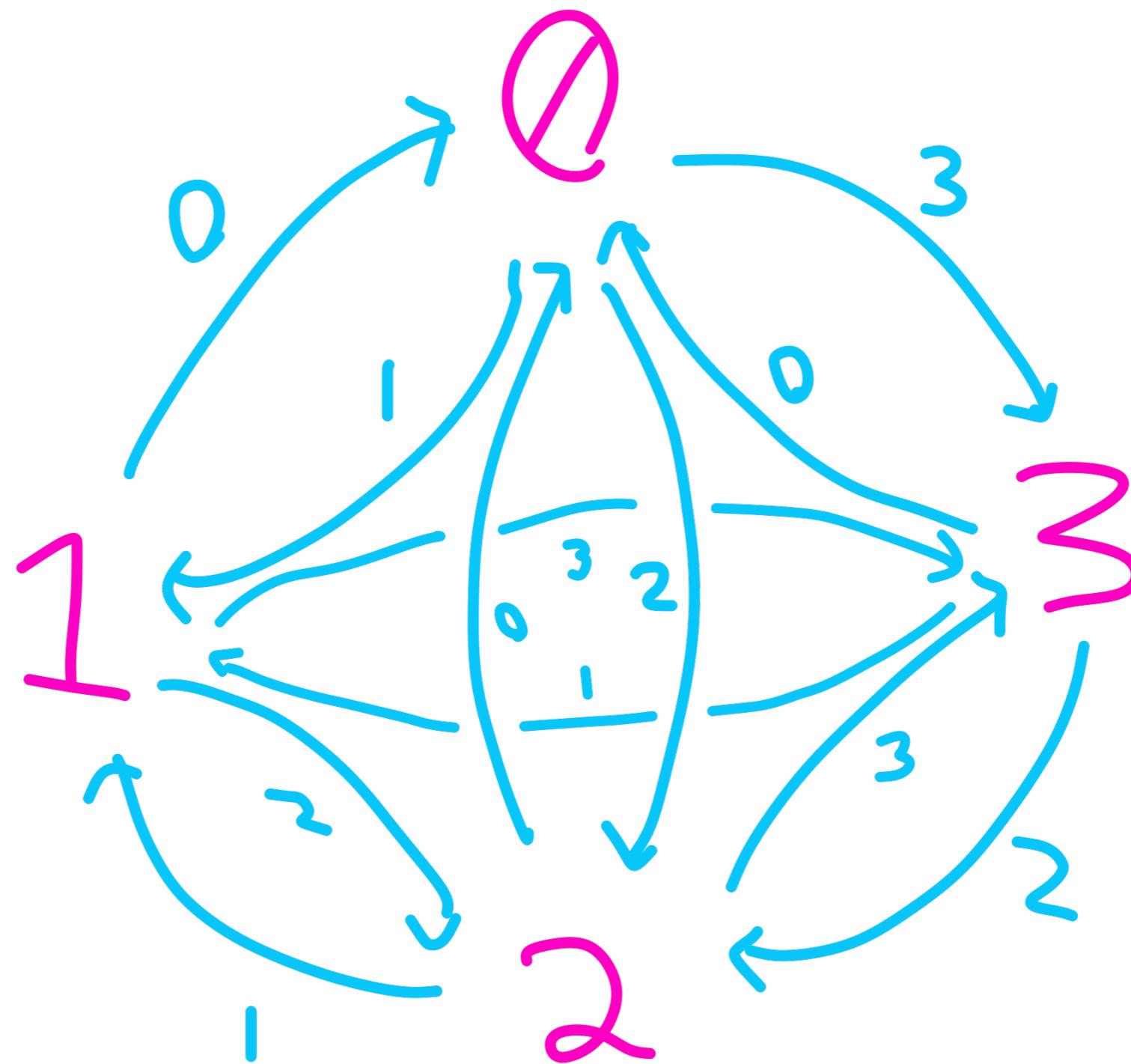
$$\boxed{5} - w(3) \rightarrow \boxed{3}$$

$$\boxed{5} - w(3) \rightarrow \boxed{8}$$

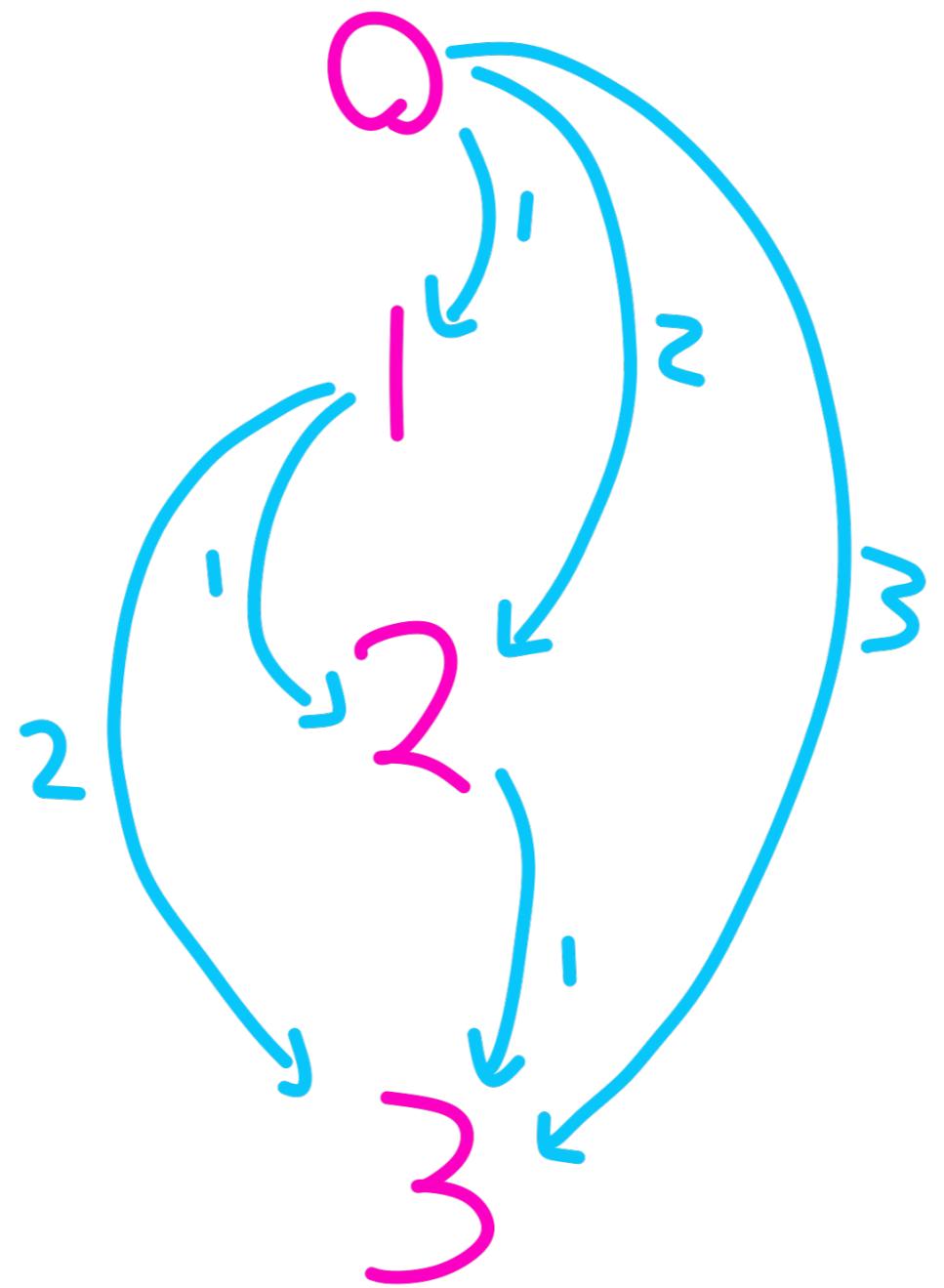
$$\boxed{\{1, 2\}} - w(3) \rightarrow \boxed{\{1, 2, 3\}}$$

$$\boxed{[1, 2]} - w(3) \rightarrow \boxed{[1, 2, 3]}$$

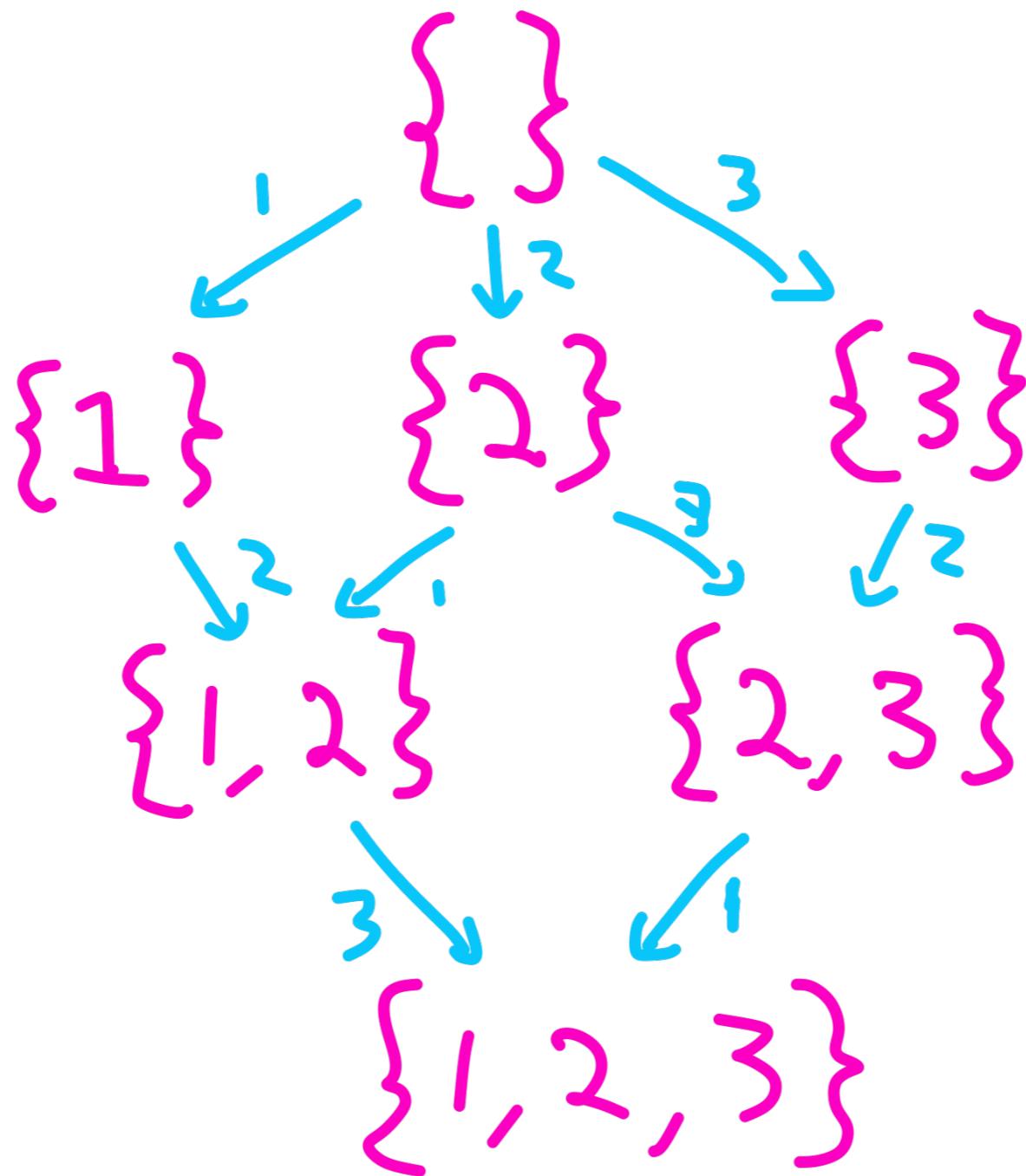
Registers



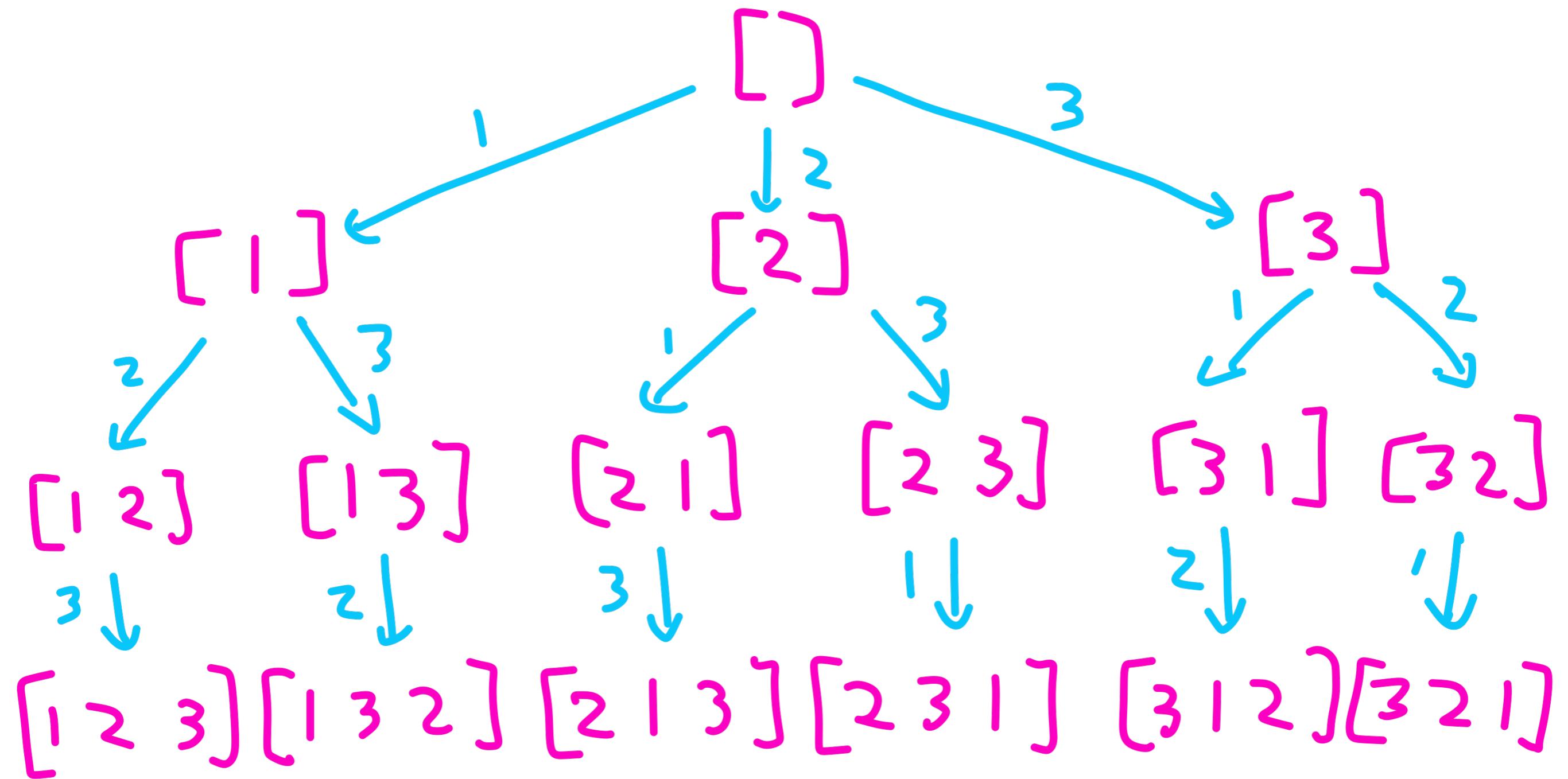
Counters

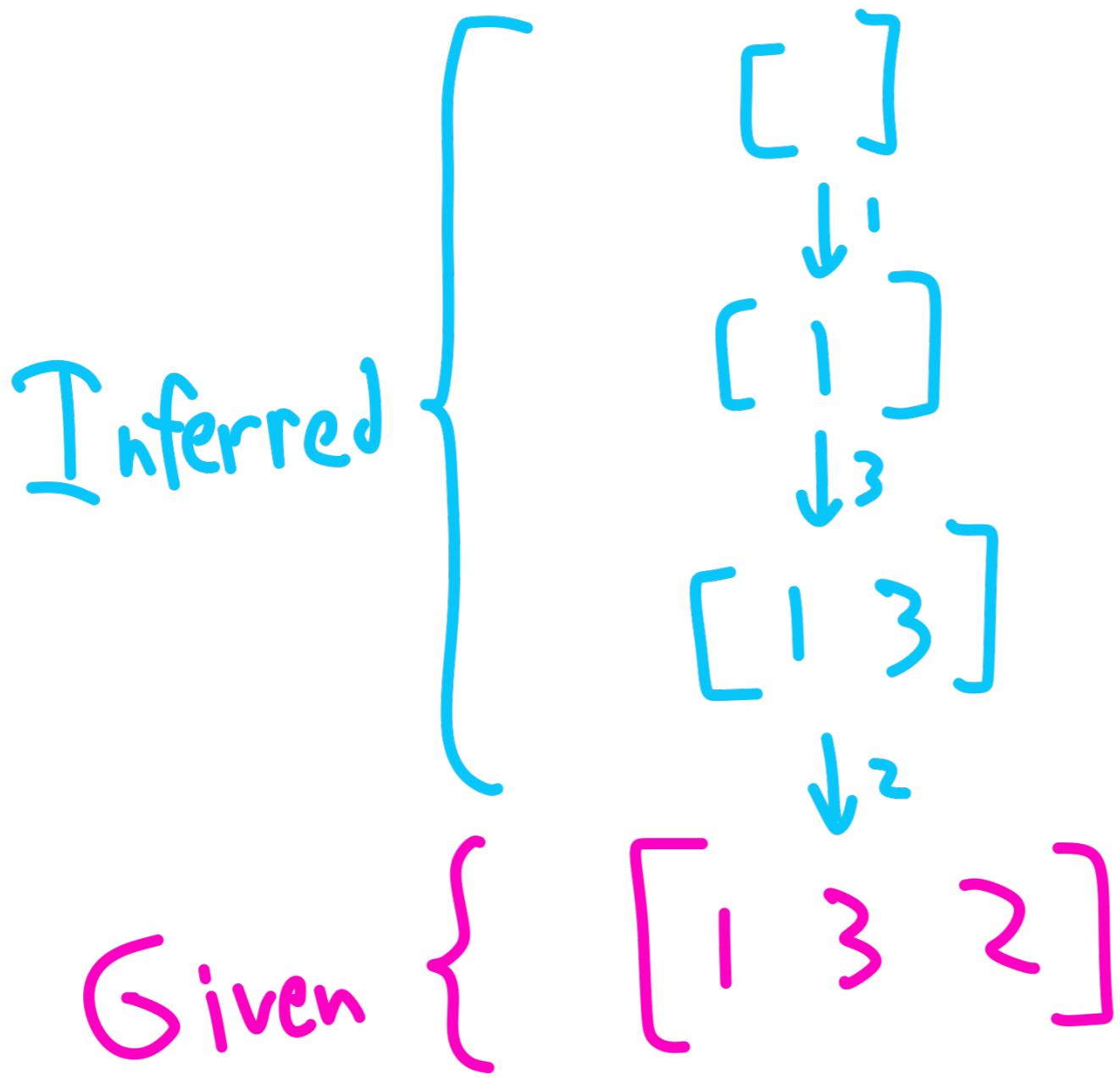


Sets



Lists





Traceability

- Lists w/append!

- Logs w/insert !

- Strings w/concat !

T+`s in the SQL standard!

EXAMPLES

$\bar{T}_1 : w(x, 1), w(y, 1)$

$T_2 : w(x, 2), w(y, 2)$

$$\begin{array}{l} \overline{T}_1 : w(x, 1), w(y, 1) \\ \downarrow \text{ww} \\ T_2 : w(x, 2), w(y, 2) \end{array} \quad \boxed{\text{GO}} \quad \boxed{\text{NO}}$$

$\bar{T}_3 : r(x, [1, 2]), r(y, [2, 1])$

- $T_1: r(x, [0, \underline{1}]), w(x, 2)$
 $\uparrow wr$
 $\downarrow ww$
]
 $G1^c$
- $\bar{T}_2: w(x, \underline{1}), w(x, 3)$
- $T_3: r(x, [0, 1, \underline{2}, \underline{3}]),$
⋮

$T_1 : r(x, [1]), w(y, 1)$

$T_2 : w(x, 2), r(y, [])$

$T_3 : r(x, [\underline{1}, \underline{2}])$

G2

Well yes, but

HONX

T_1

T_2

T_3

T_{11}

T_8

T_4

T_{14}

T_{12}

T_7

T_{13}

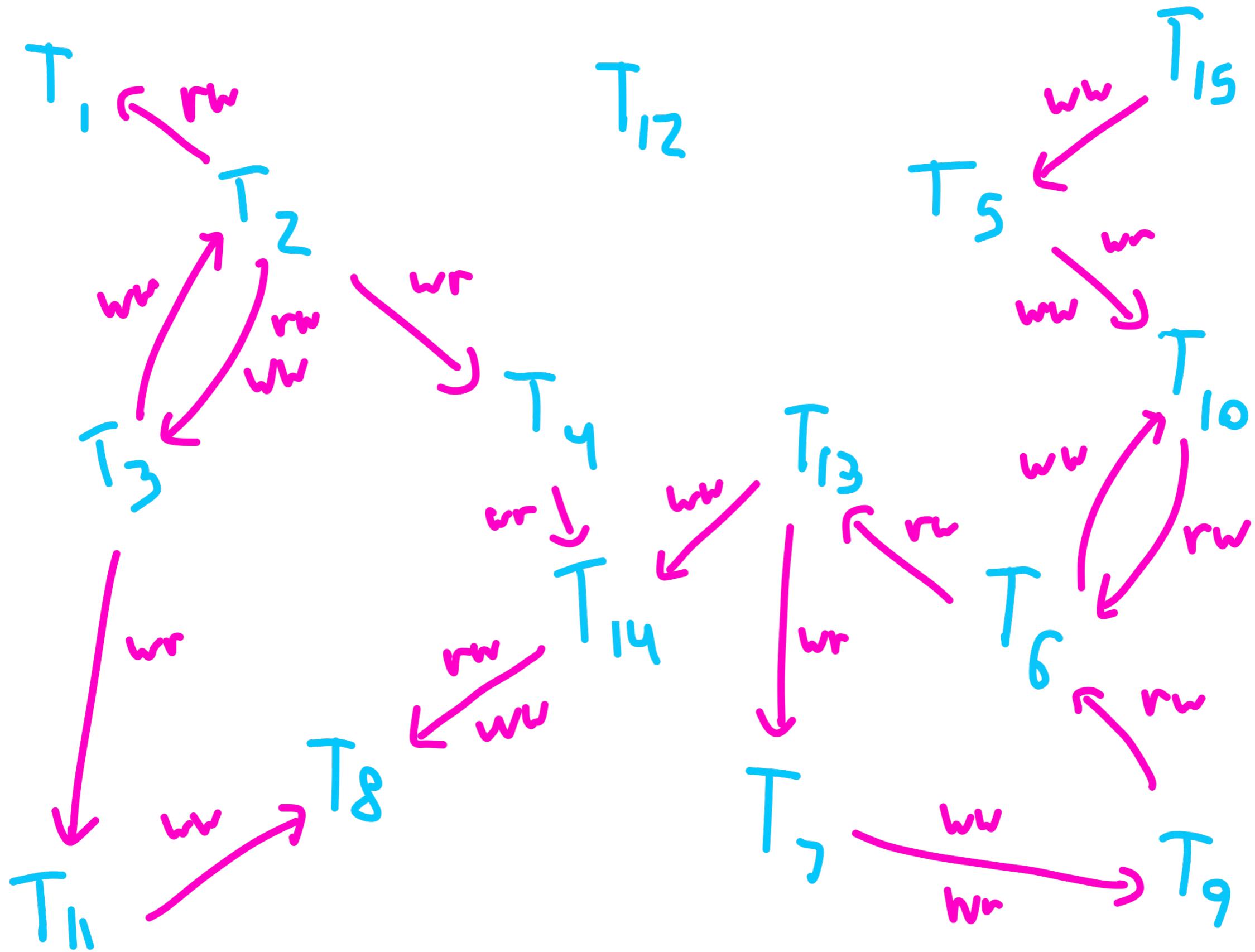
T_5

T_9

T_{10}

T_6

T_{15}

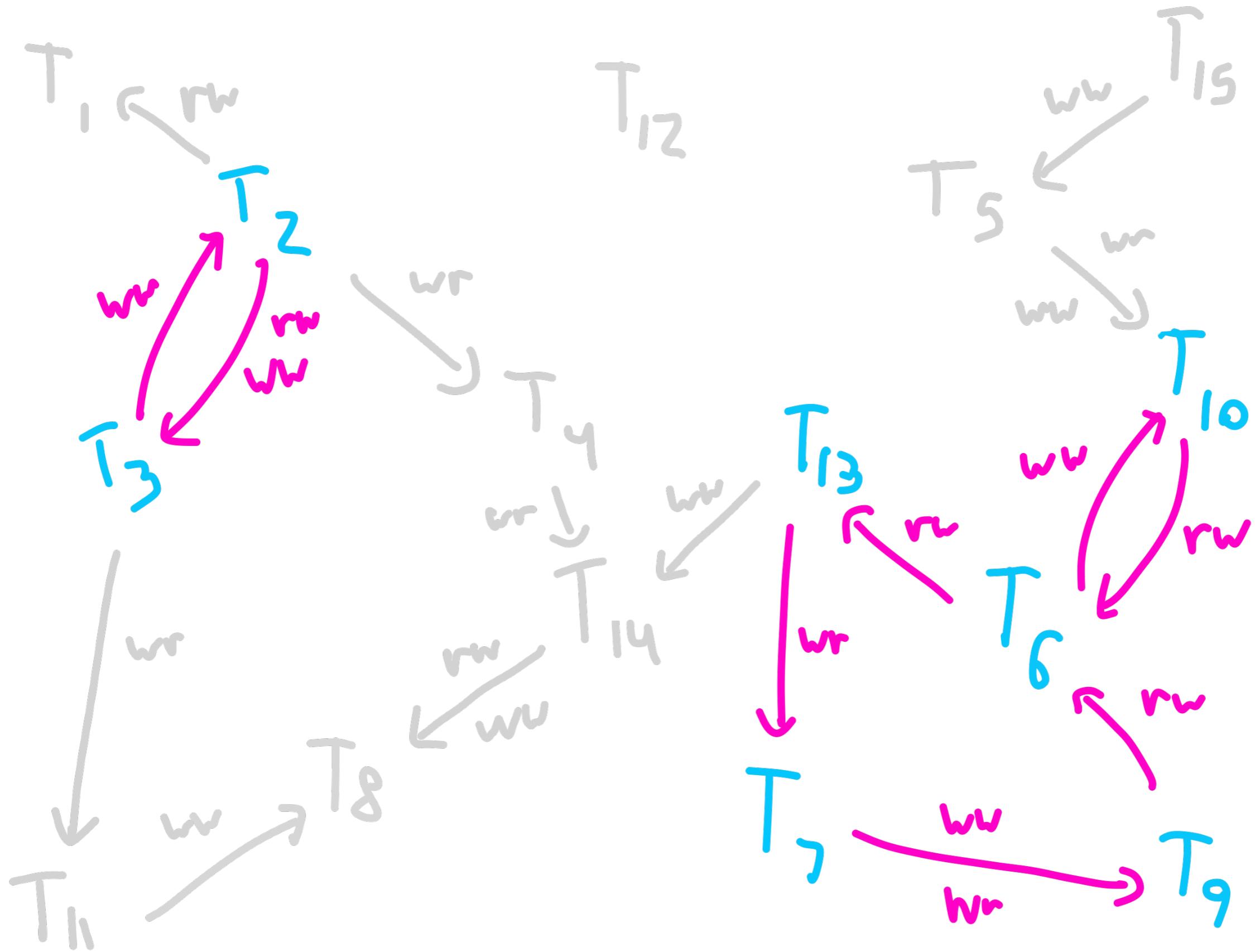


Tarjan 1972

Strongly Connected

Components in

$O(V + E)$

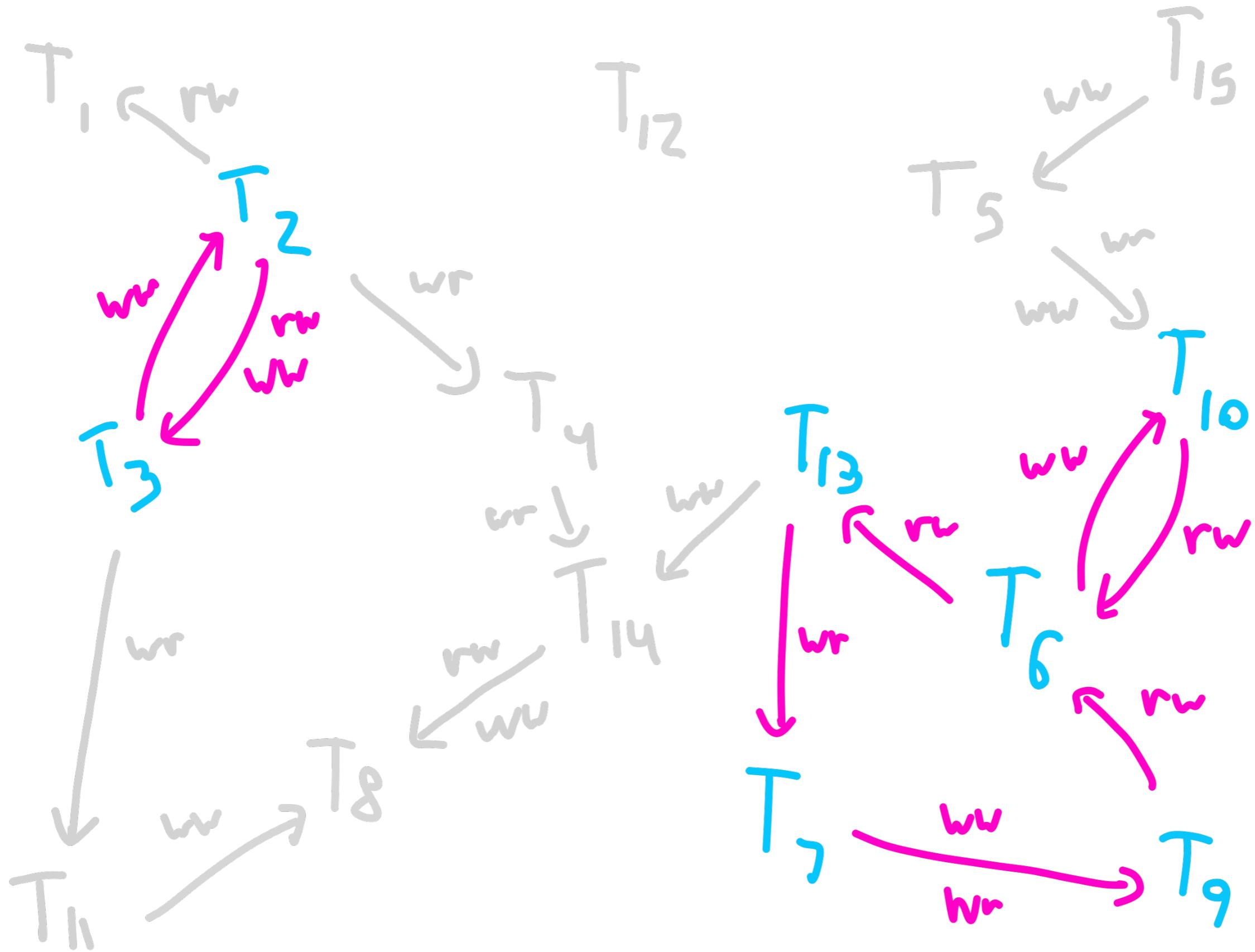




$T_2 \xrightarrow{ww} \bar{T}_3 \xrightarrow{ww} \bar{T}_2 : G\emptyset$

$\bar{T}_2 \xrightarrow{rw} T_3 \xrightarrow{ww} T_2 : G\text{-sing/c}$

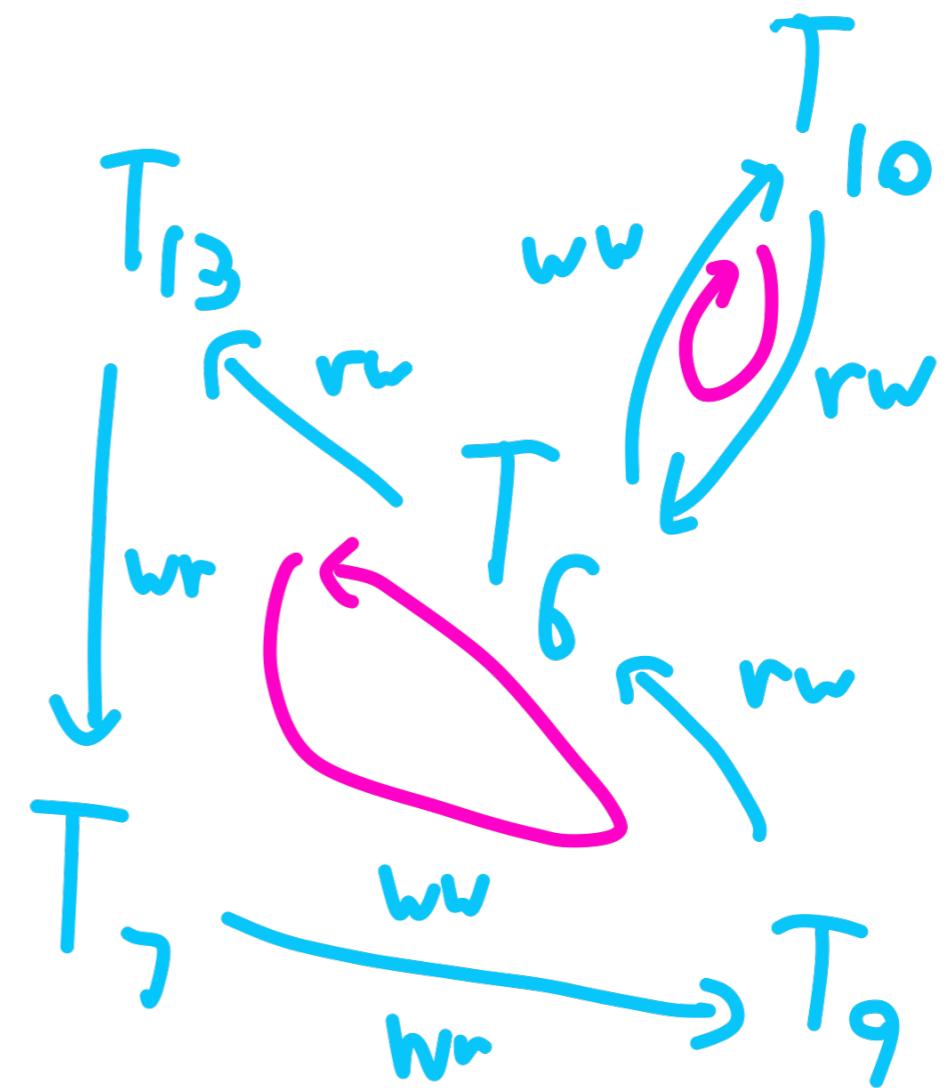
BFS



$T_{10} \xrightarrow{rw} T_6 \xrightarrow{ww} T_{10}$; G-Single

$T_{13} \xrightarrow{wr} T_7 \xrightarrow{w} T_9 \xrightarrow{rv} T_6 \xrightarrow{rw} T_{13}$; G2

BFS



Automatic
Explanations

$T_1 < T_2$ because T_2 observed T_1 's
write of 3 to x.

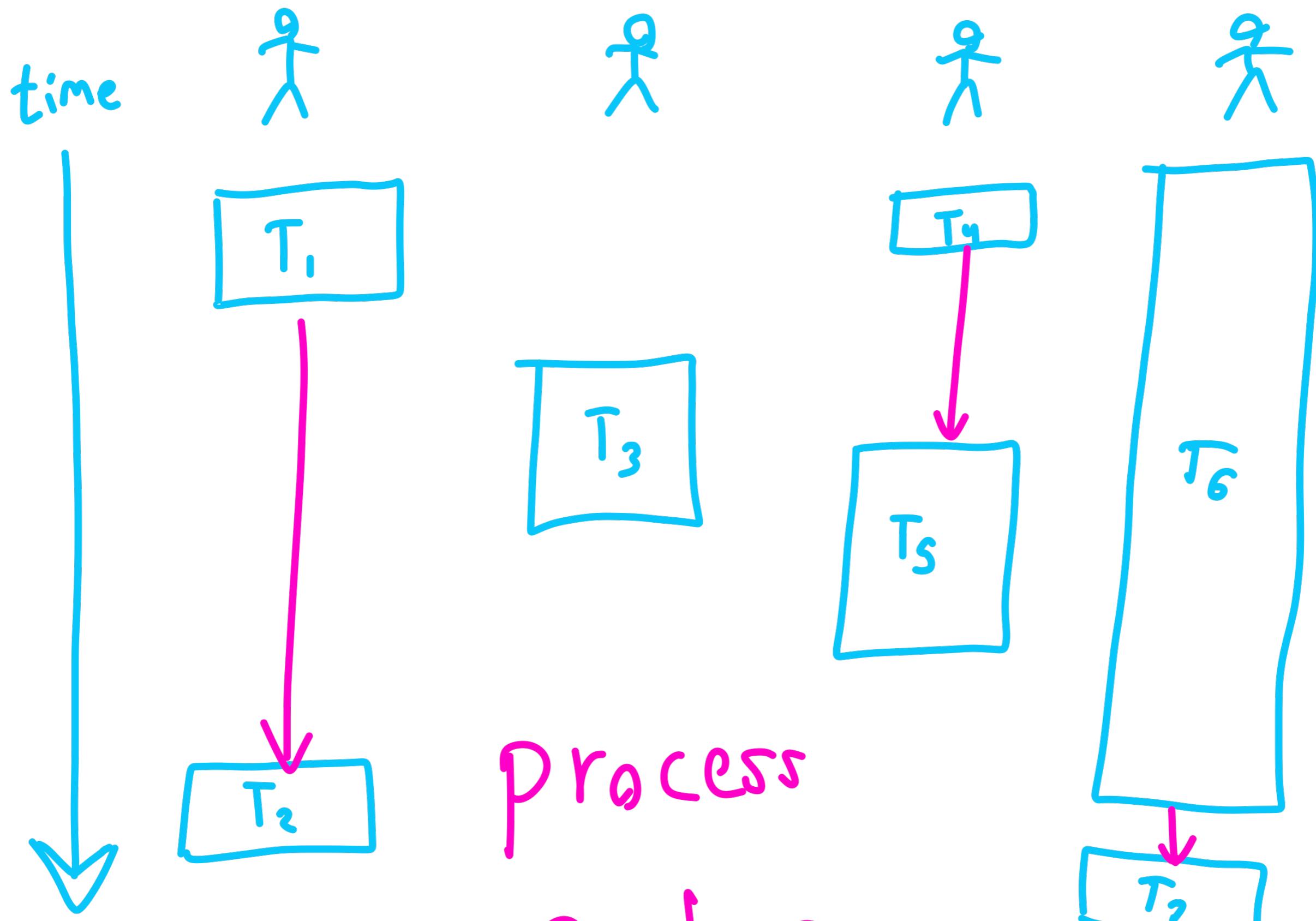
$T_2 < T_3$ because T_2 did not observe
 T_3 's write of 5 to y.

$T_3 < T_1$ because T_1 appended 3 to
z before T_3 appended 7.

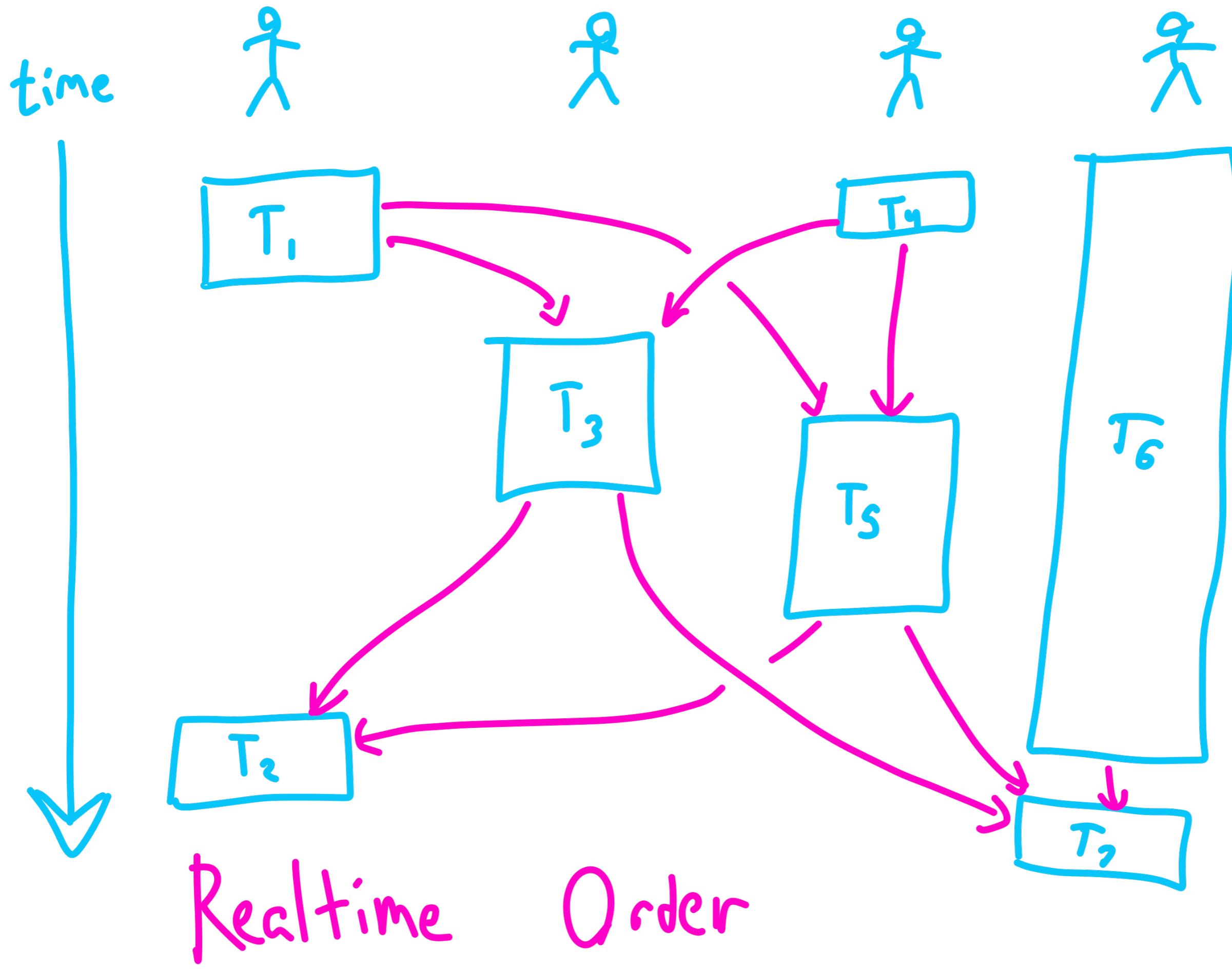
A contradiction! \square

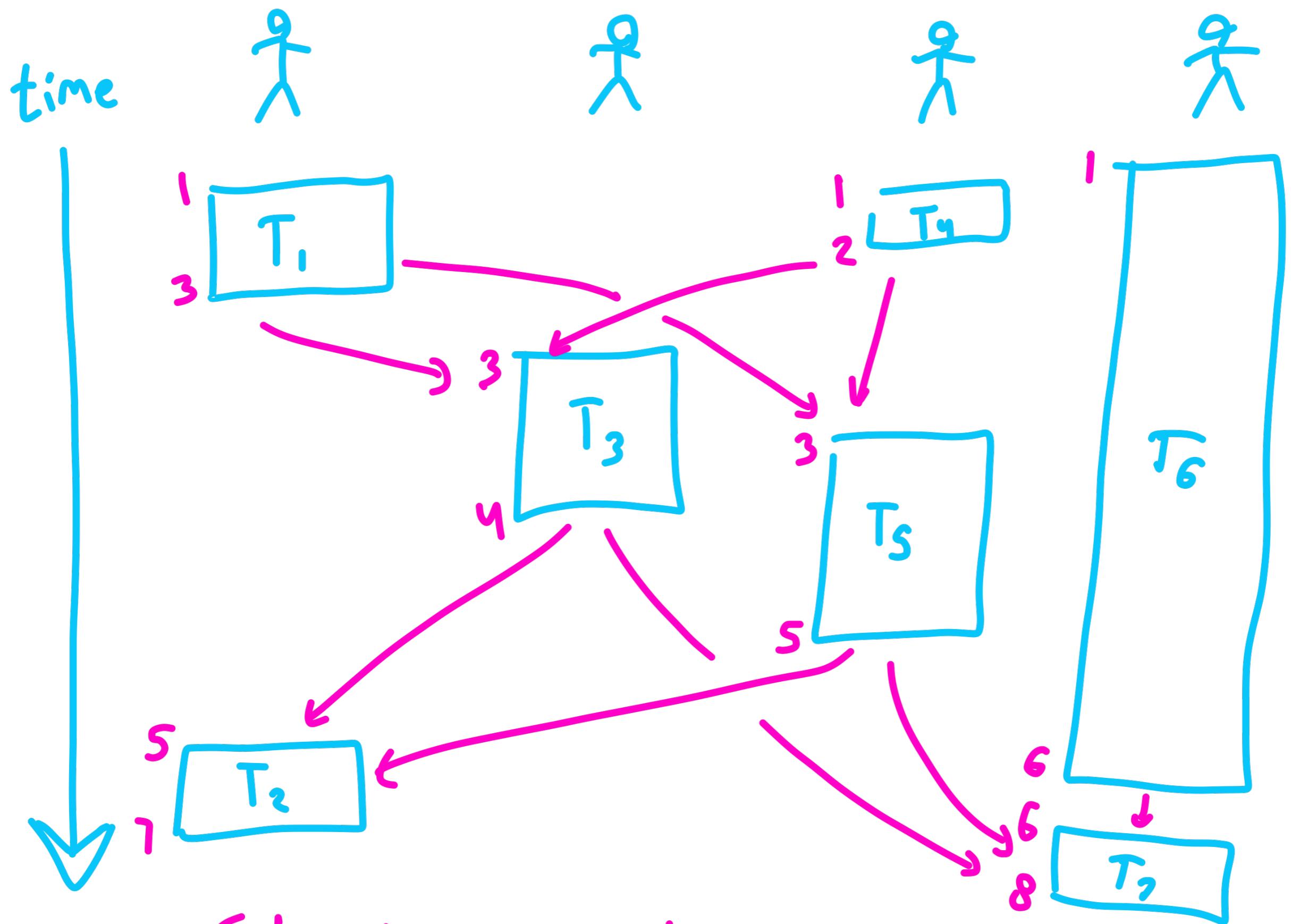
- more -

C O N S T R A I N T



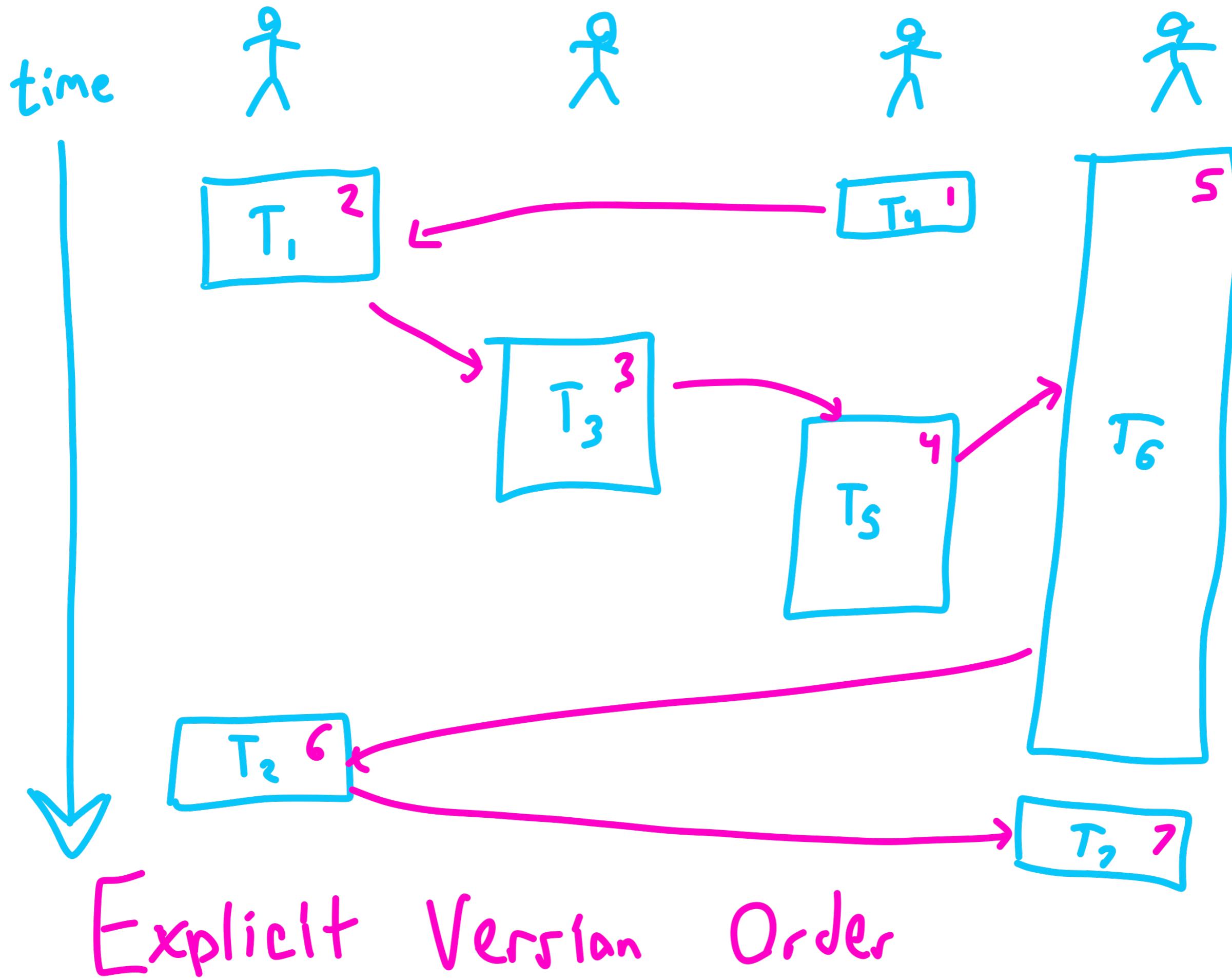
Process
Order





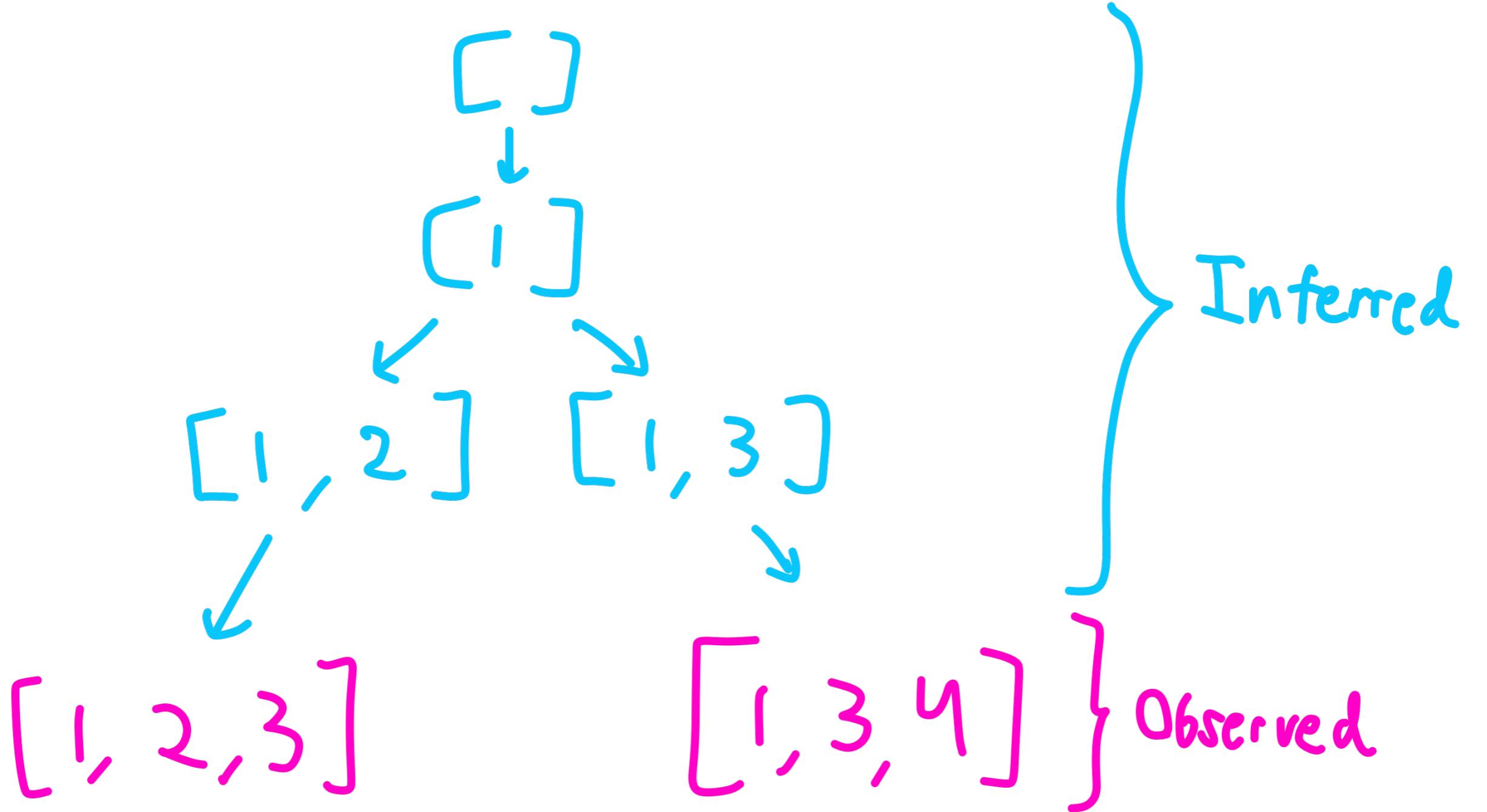
Start - Ordered

Serialization Graph



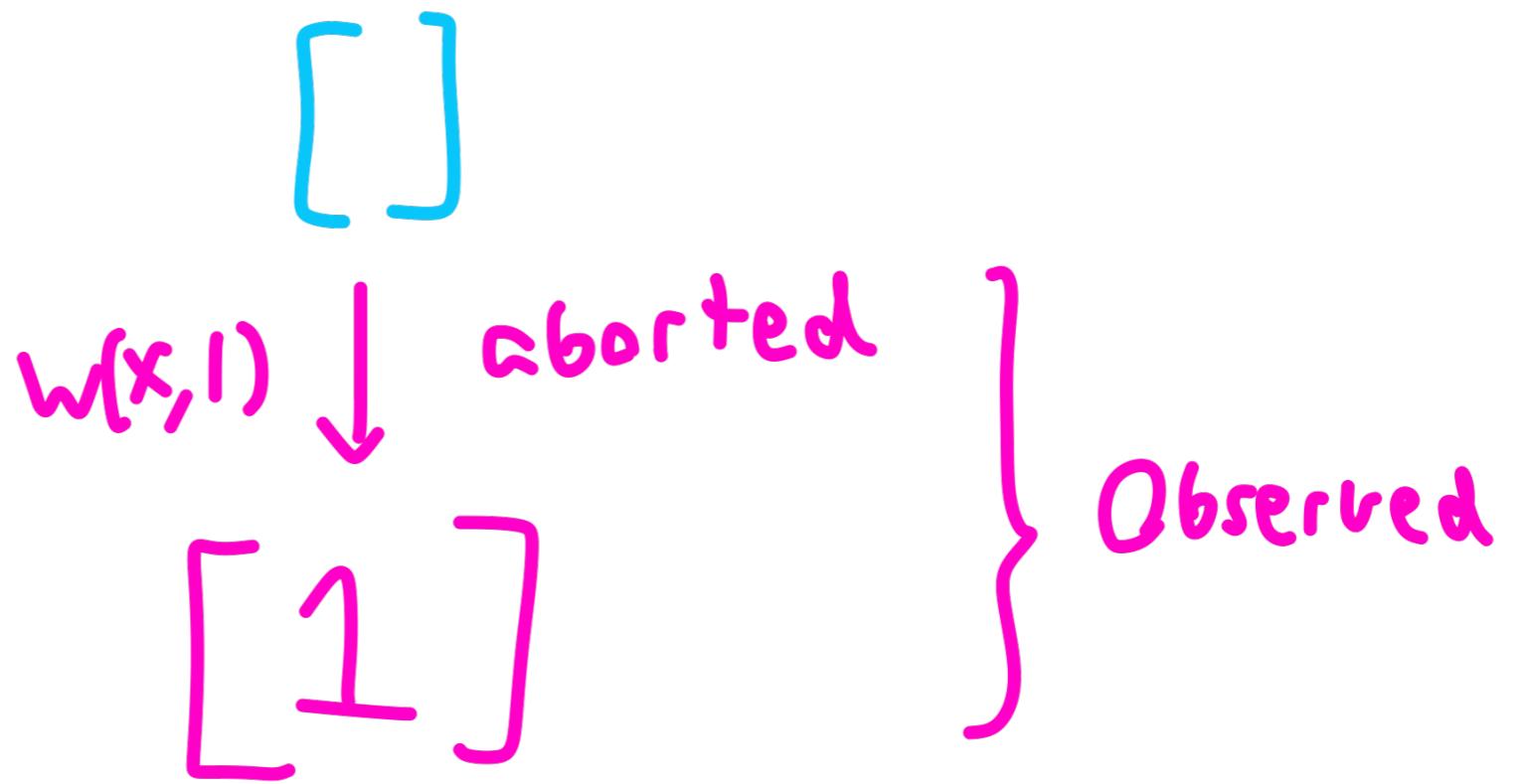
Other

ANOMALIES

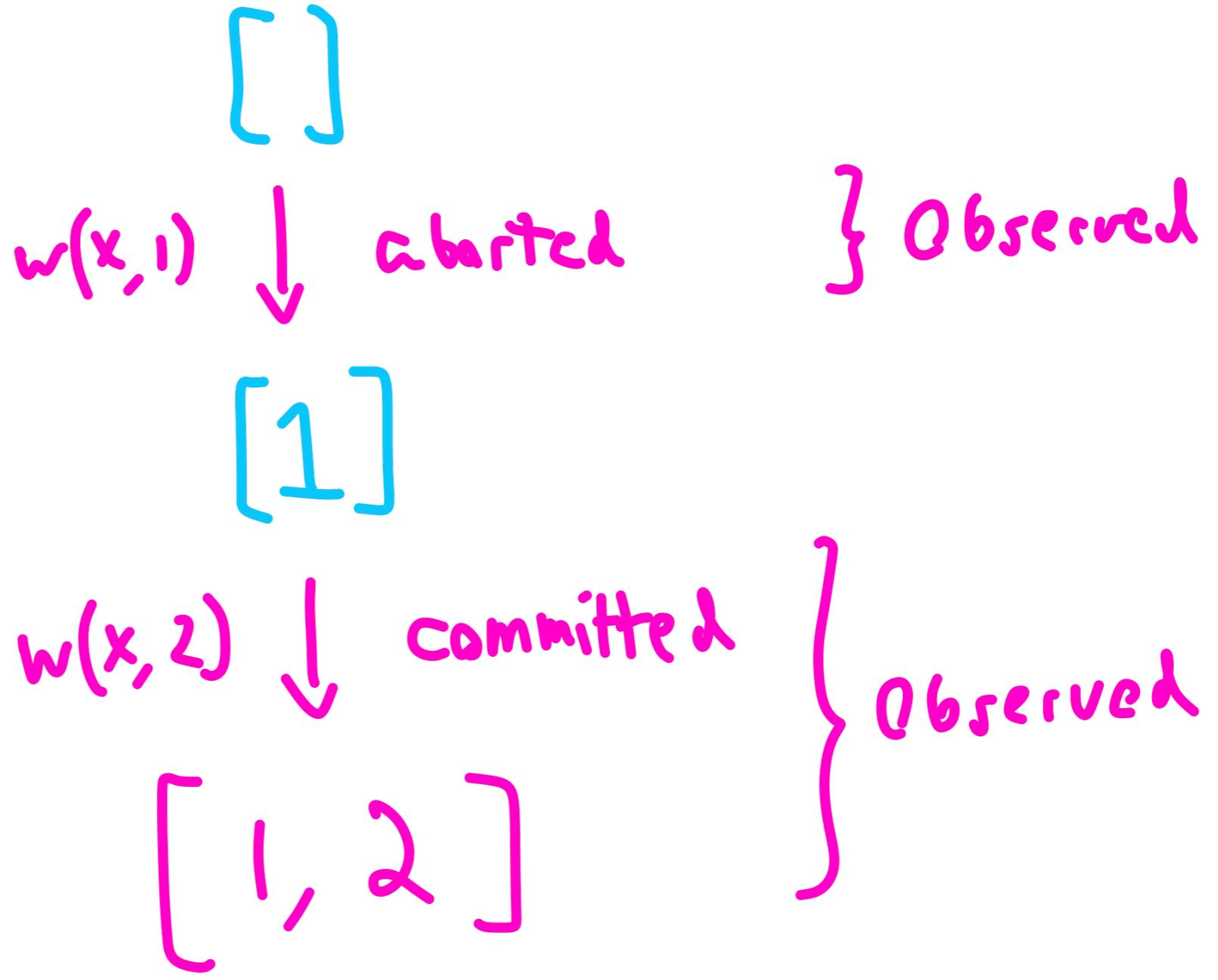


Aborted

Read



Aborted Read



Dirty Update

[1, 1, 2]

Duplicate Element

[1, ~~00~~, 3]

Garbage Element

$T_i : w(x, 3), \dots, r(x, [\dots, 2])$

Internal Inconsistency

Does it

WORK?

TiDB

2.1.7 - 3.0.0-beta.1

- Claim: SI

- Found : G2

G-single

Lost Updates

Aborted Reads

TiDB

2.1.7 - 3.0.0-beta.1

- Claim: Select... for update

⇒ No write skew

- Found: Write skew on

new rows

Yugabyte DB 1.3.1

- Claim: Serializable
- Fowl: G2-item when master nodes paused/crashed

Fauna DB 2.6.0

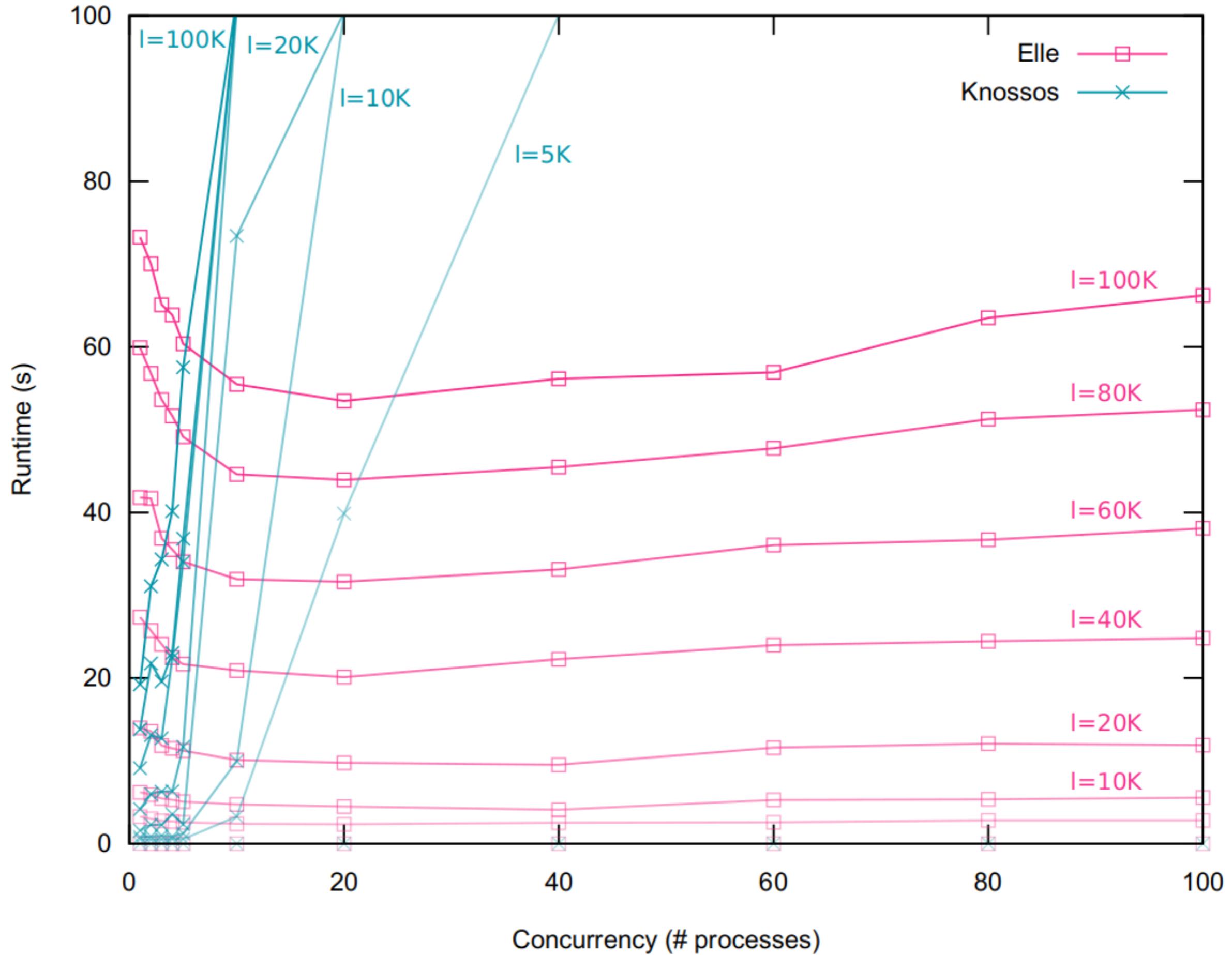
- Claim: SI \rightarrow Strict-1SR
- Found: Internal Inconsistencies in
 - index reads; did not
 - reflect txn's prev writes

And it's

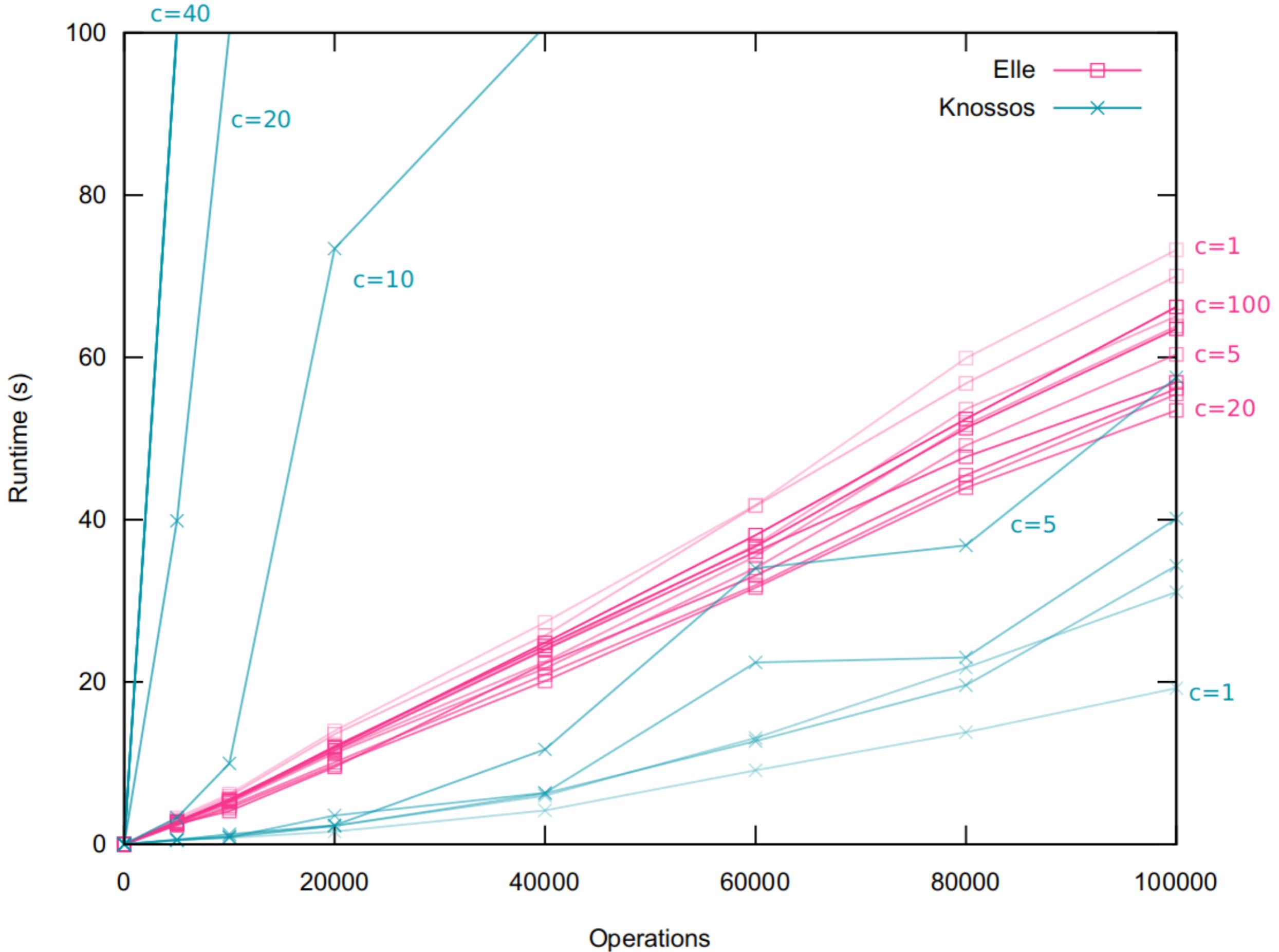
Fast

too!

Runtime vs concurrency, for various lengths (l=0, 5K, ...)

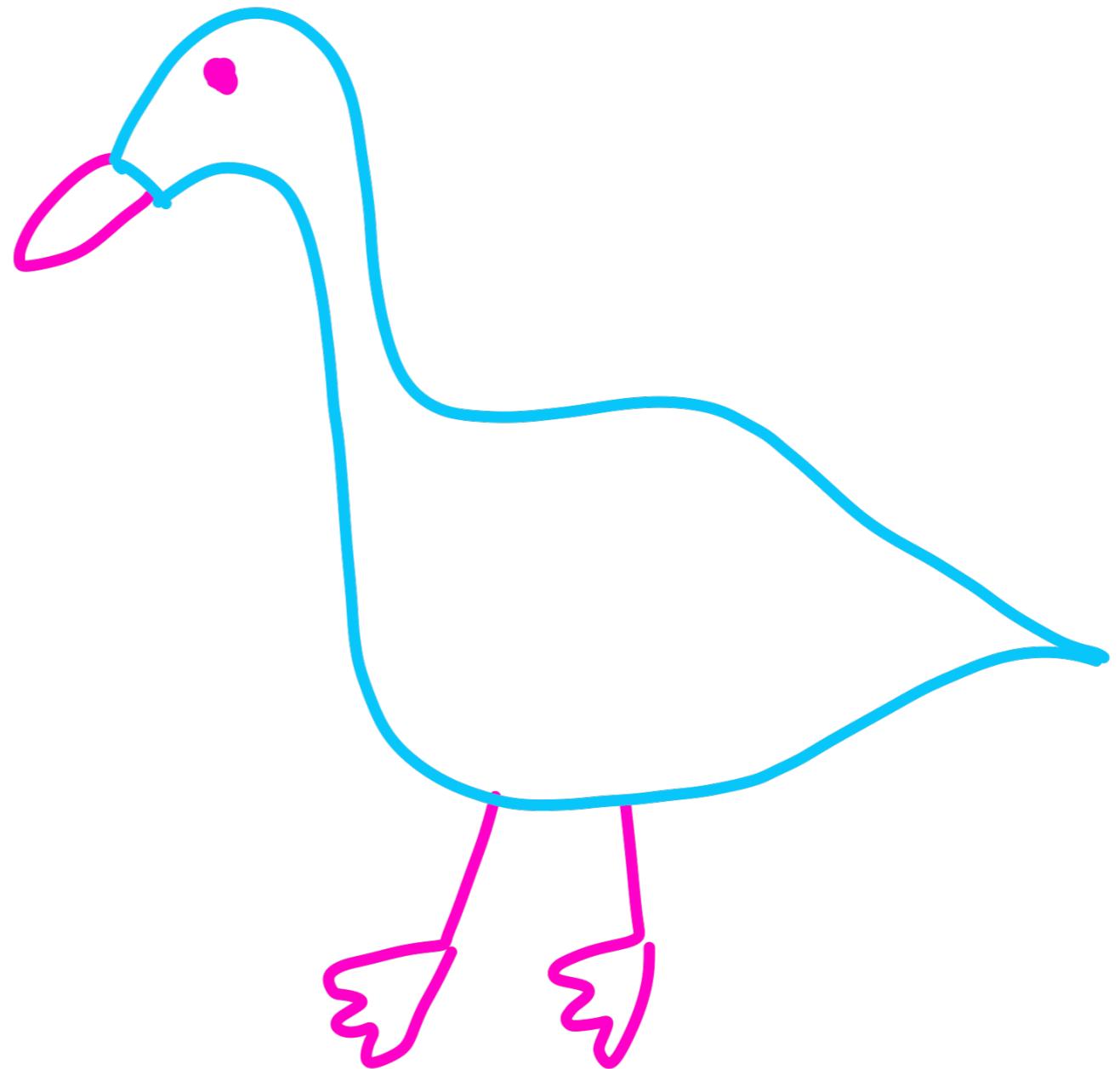


Runtime vs history length, for various concurrencies ($c=1, 2, \dots$)



ELLE

- General
- Effective
- Efficient
- Sound
- Explainable



Thanks

Peter Alvaro

Asha Karim

Kit Patella