

# Relaxed Virtual Memory

## — Armv8 edition —

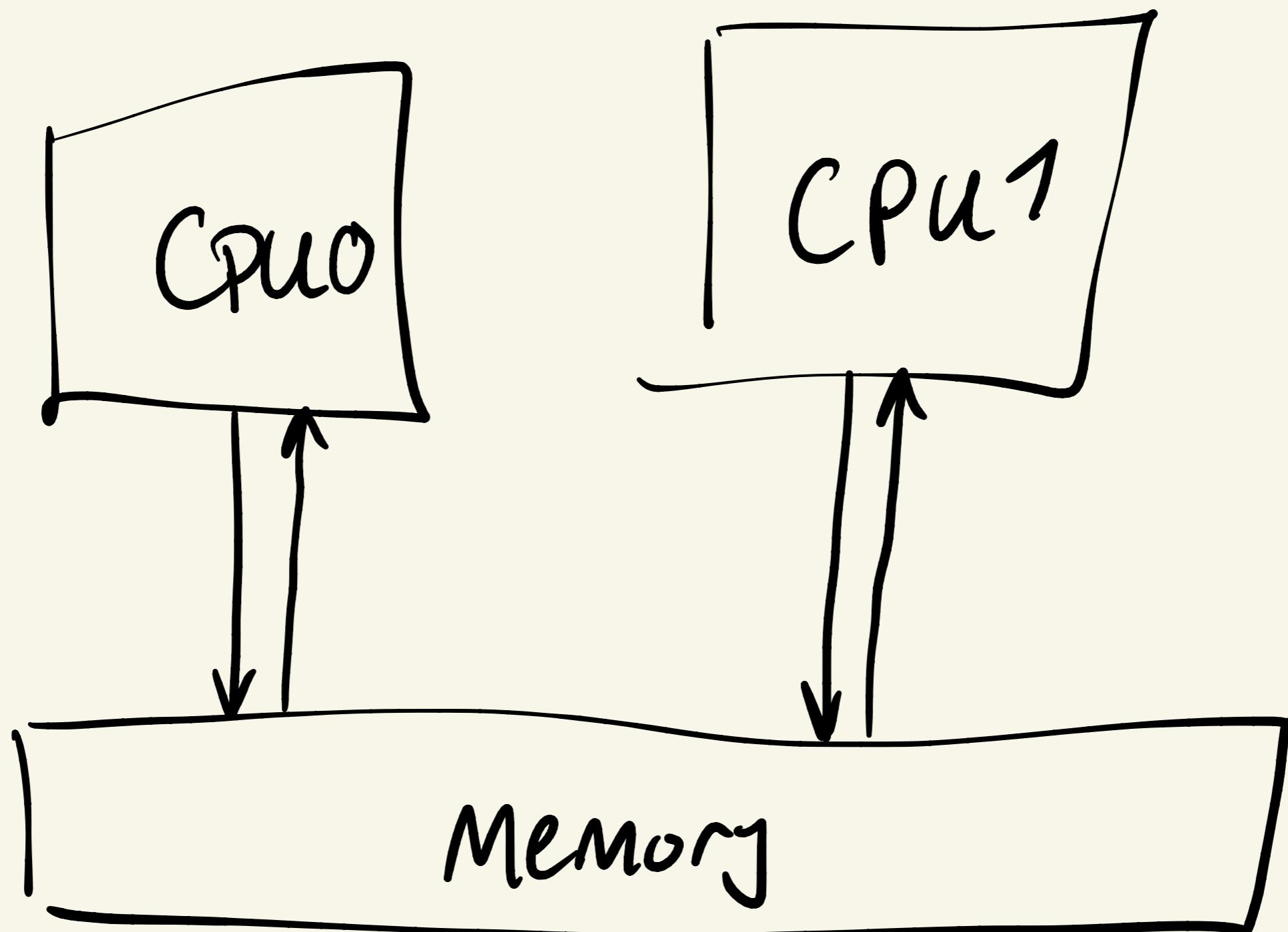
Ben Simner<sup>1</sup>, Alasdair Armstrong<sup>1</sup>, Jean Pichon-Pharabod<sup>2</sup>

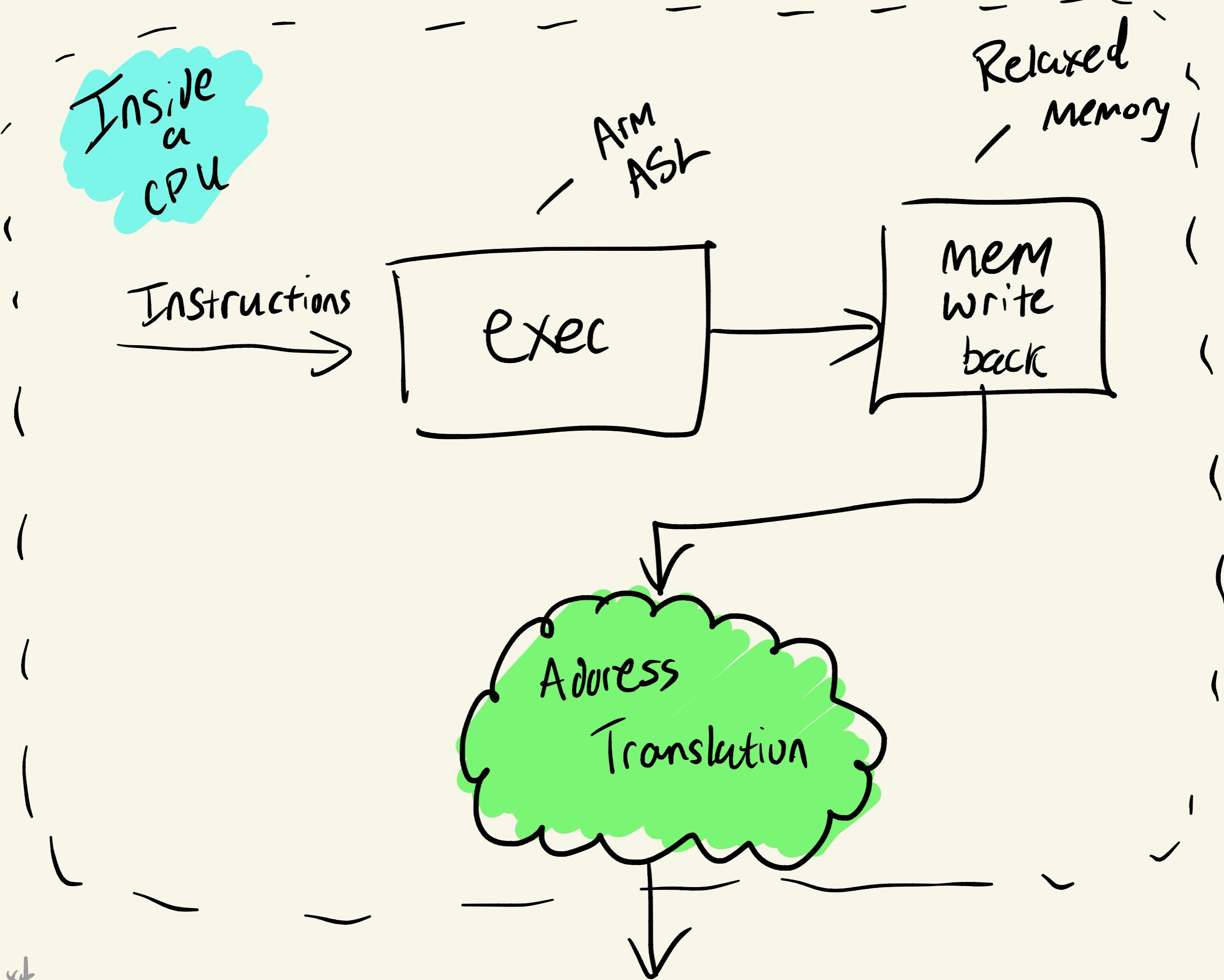
Christopher Pulte<sup>1</sup>, Richard Grisenthwaite<sup>3</sup>, Peter Sewell<sup>1</sup>

1. Cambridge, UK

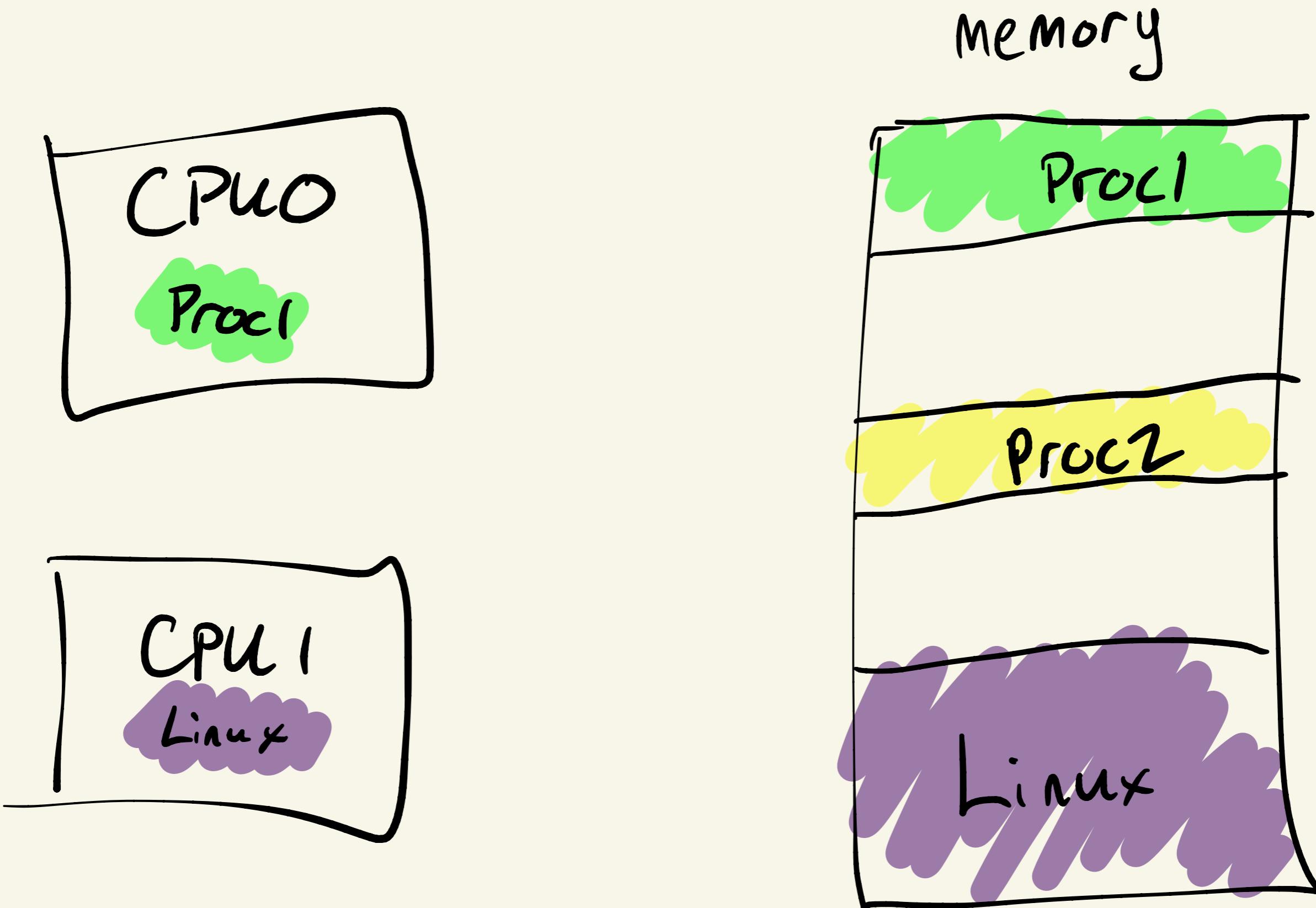
2. Aarhus, DK

3. Arm Ltd., UK

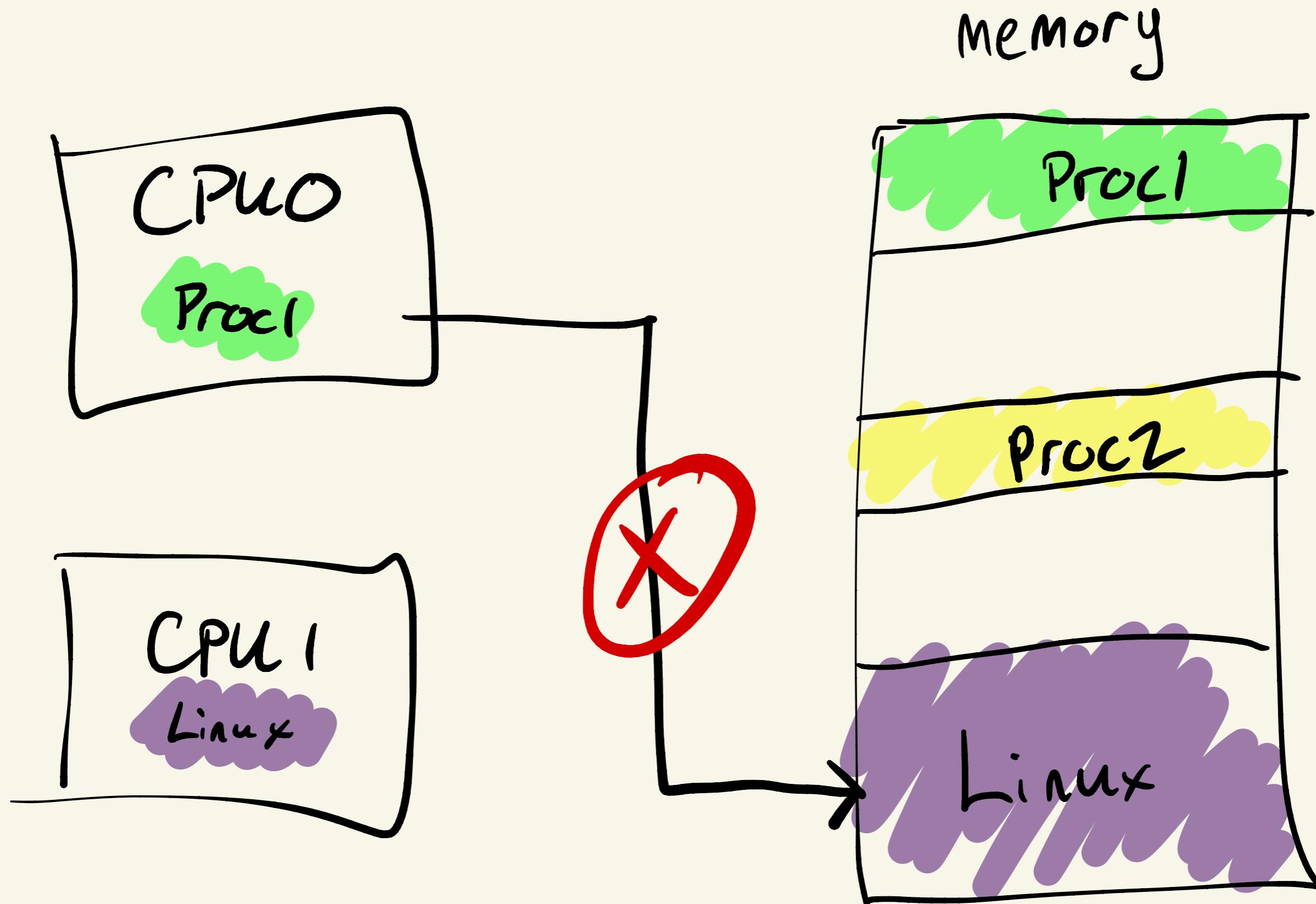




# Why Virtual Memory

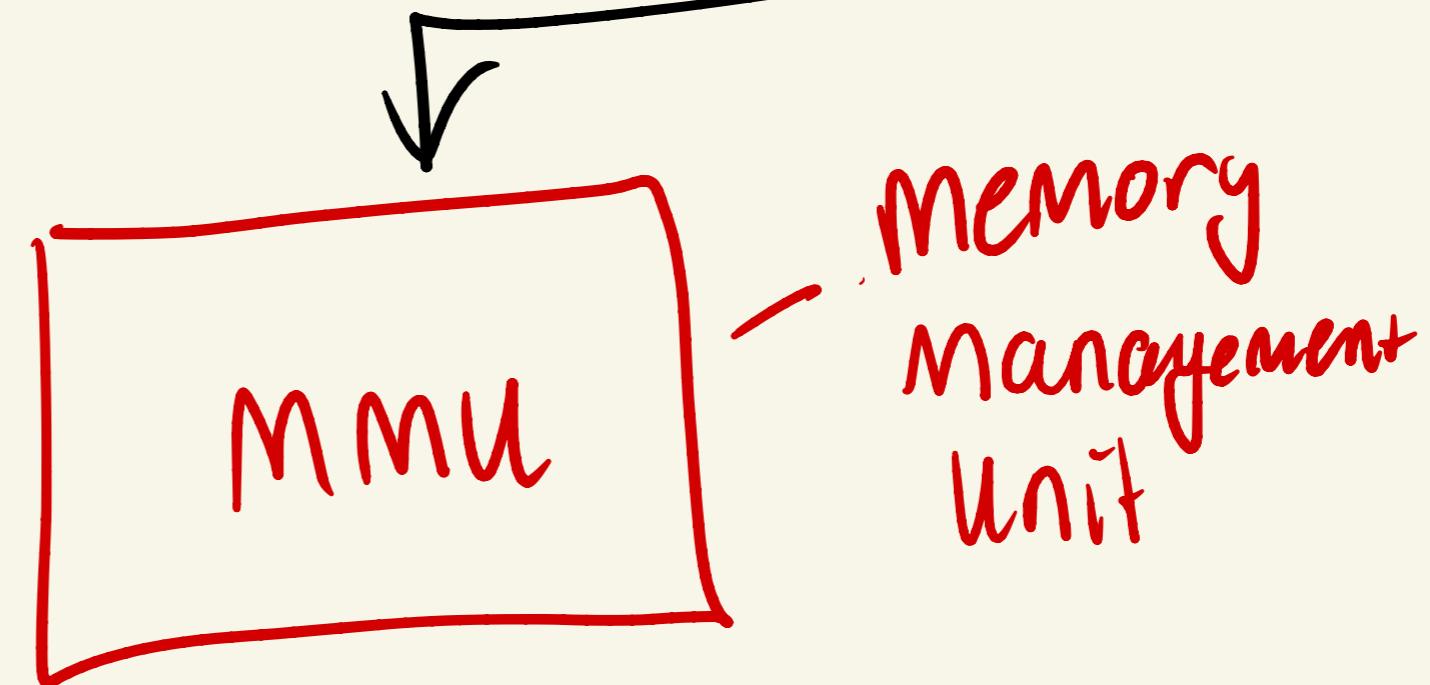
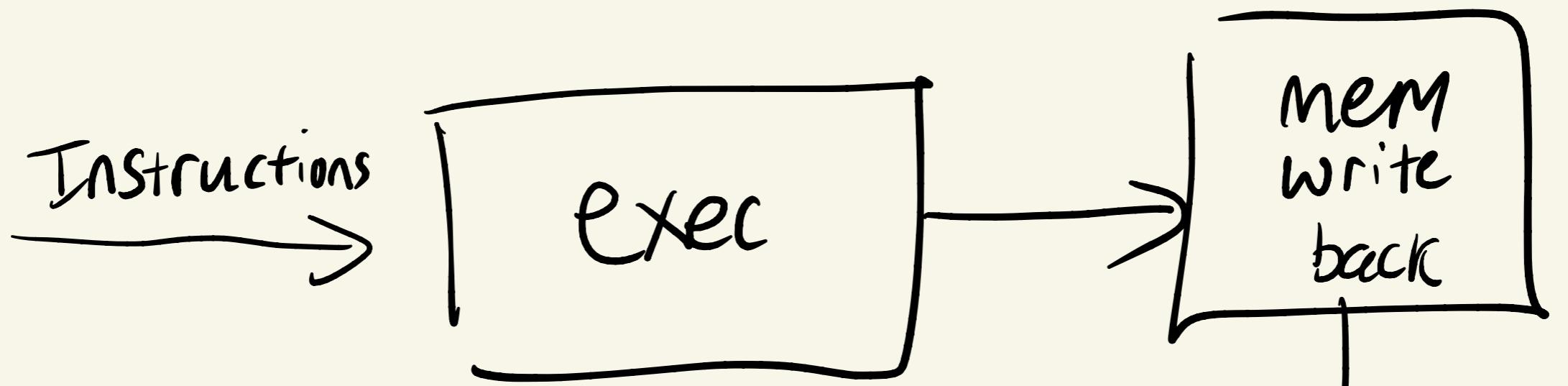


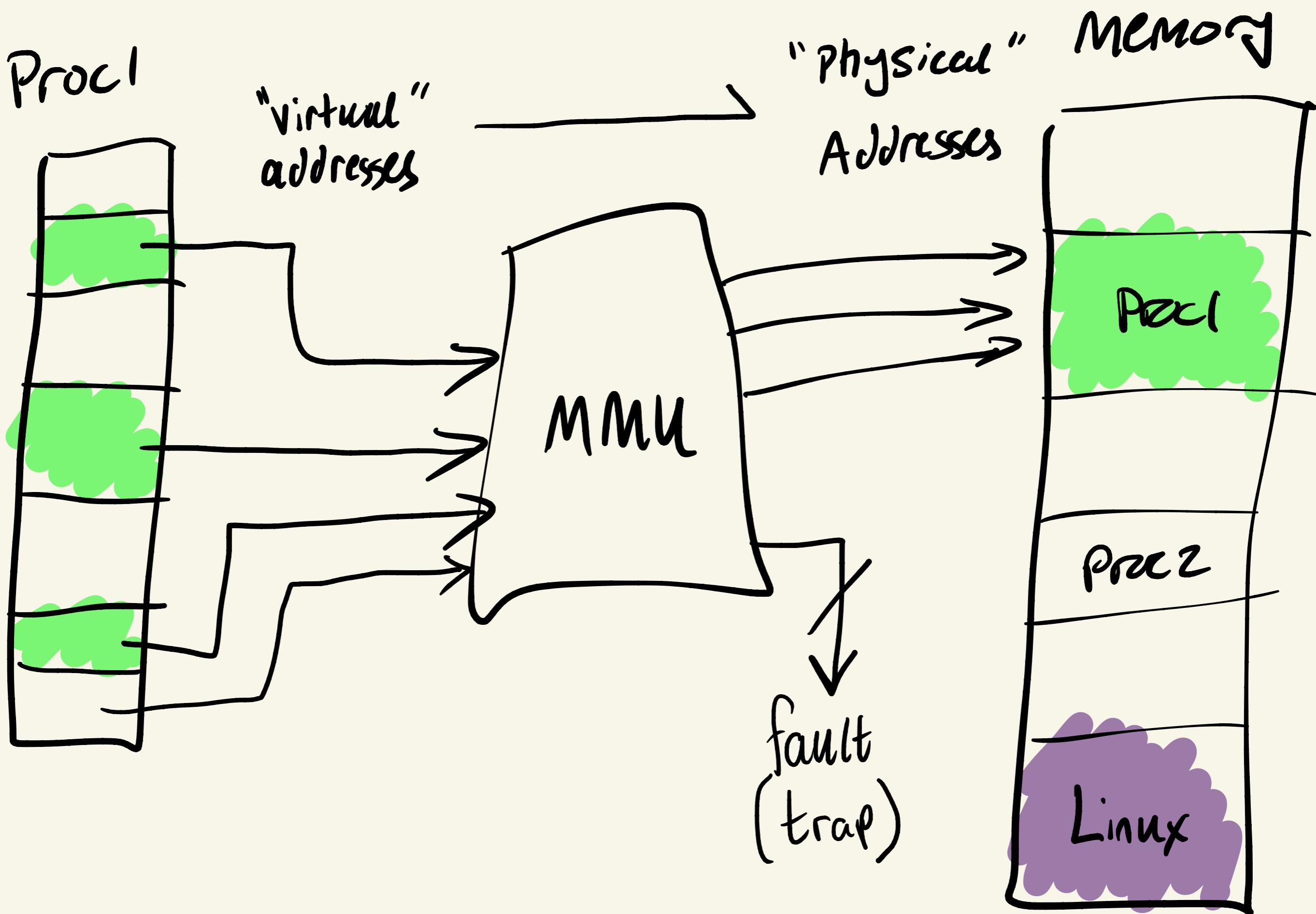
# Why Virtual Memory

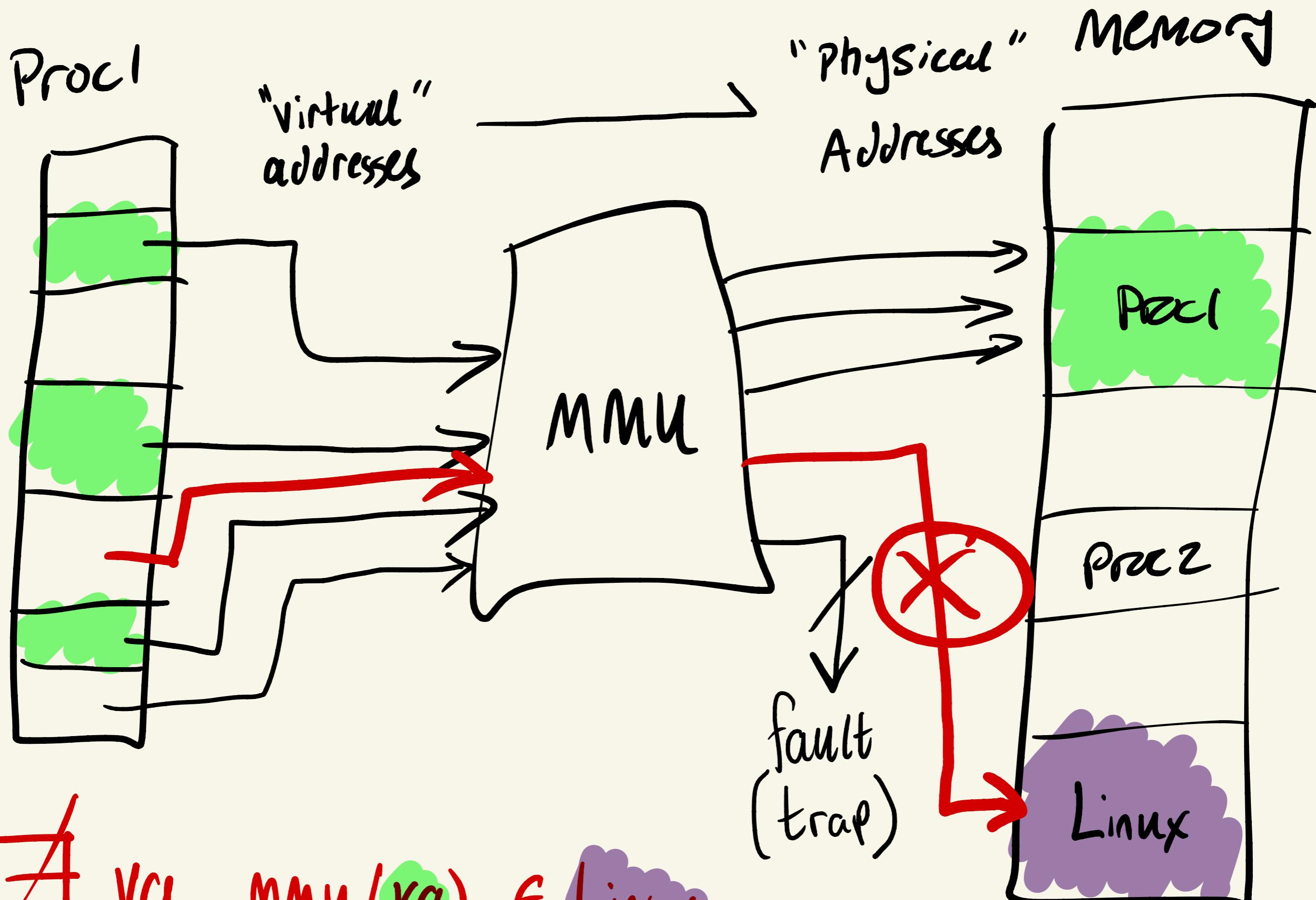


# Aims

- Formal Semantics
- Build Support tooling
- reason about real Systems code implementing isolation.



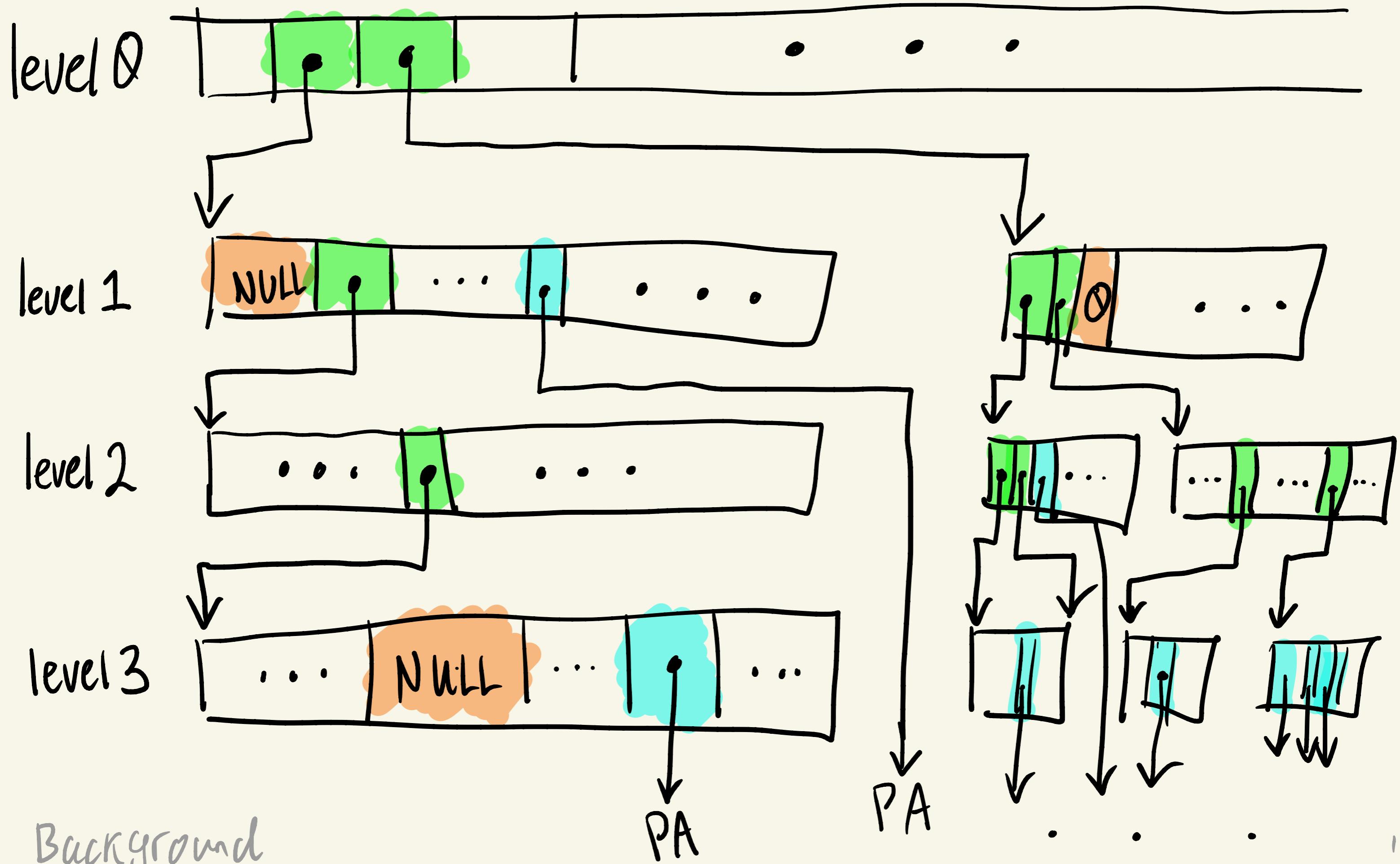




~~A~~ Va. MMU(ra) ∈ Linux

Background

# In-memory Data Structure



LDR  $r_0$ , [VA]

VA  $\rightarrow$  R PA

VA

LDR  $r_0$ , [VA]



$T_{SIL0} \dashrightarrow T_{SIL1} \dashrightarrow T_{SIL2} \dashrightarrow T_{SIL3} \dashrightarrow R$  PA

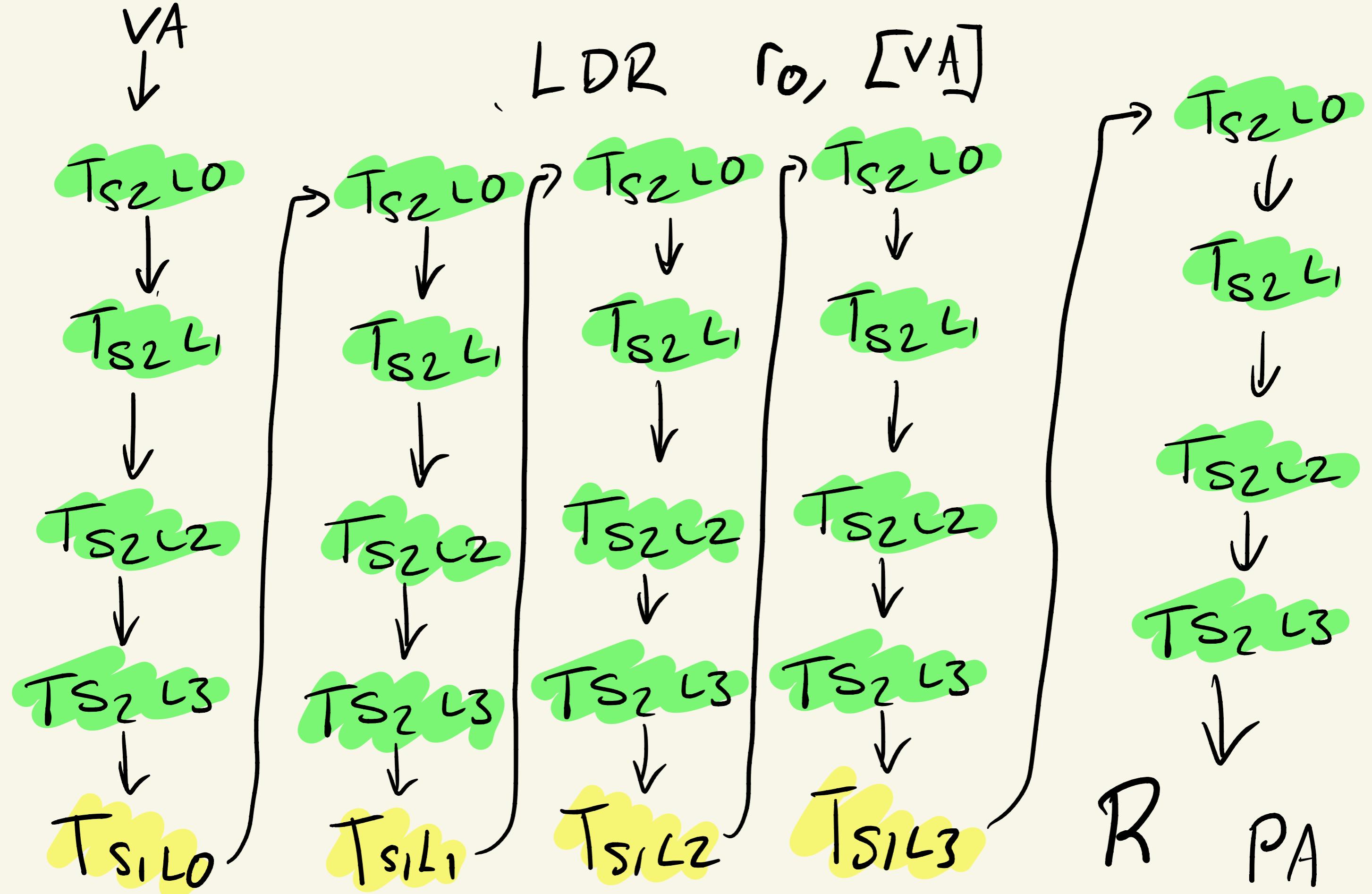
VA

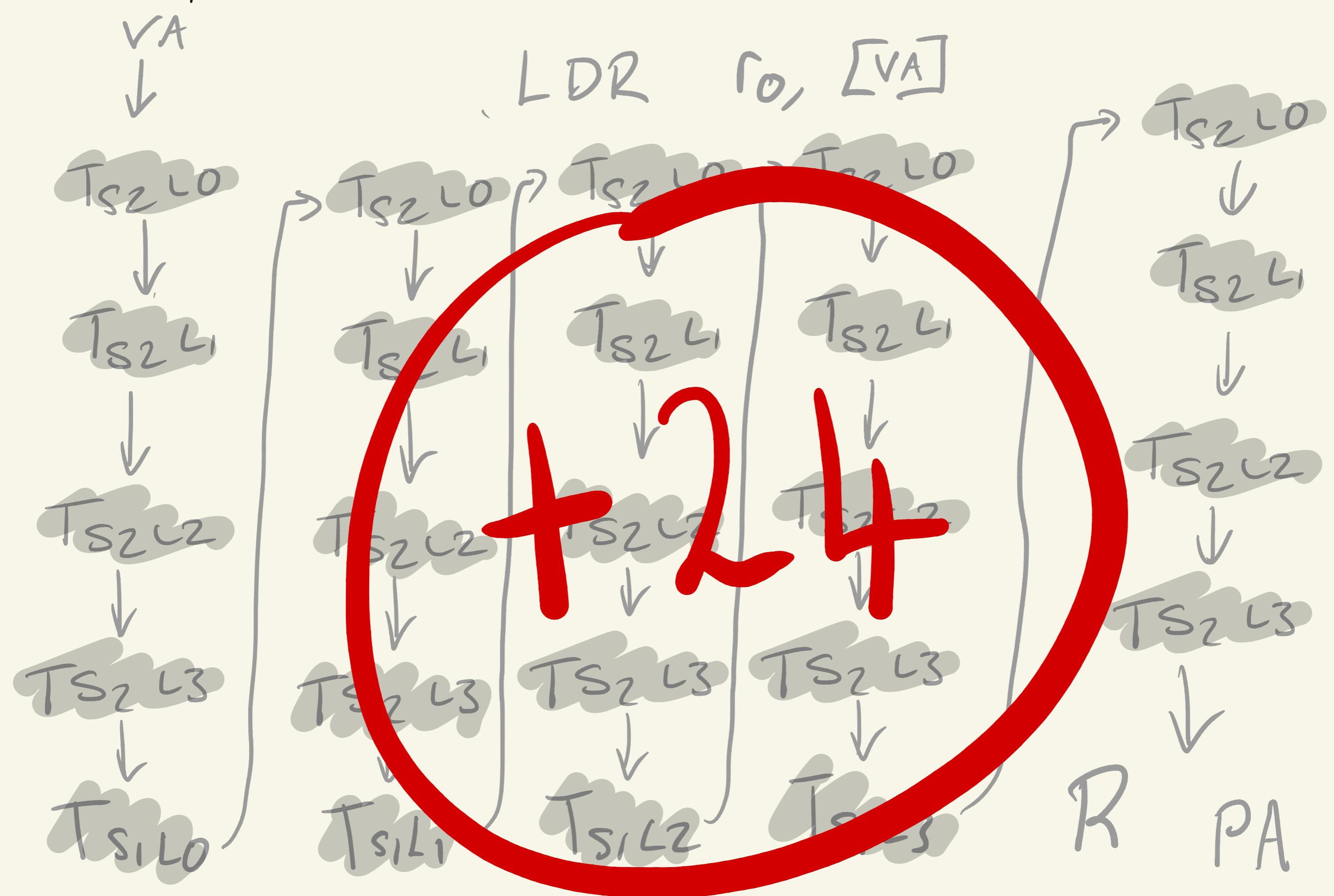
LDR  $r_0$ , [VA]

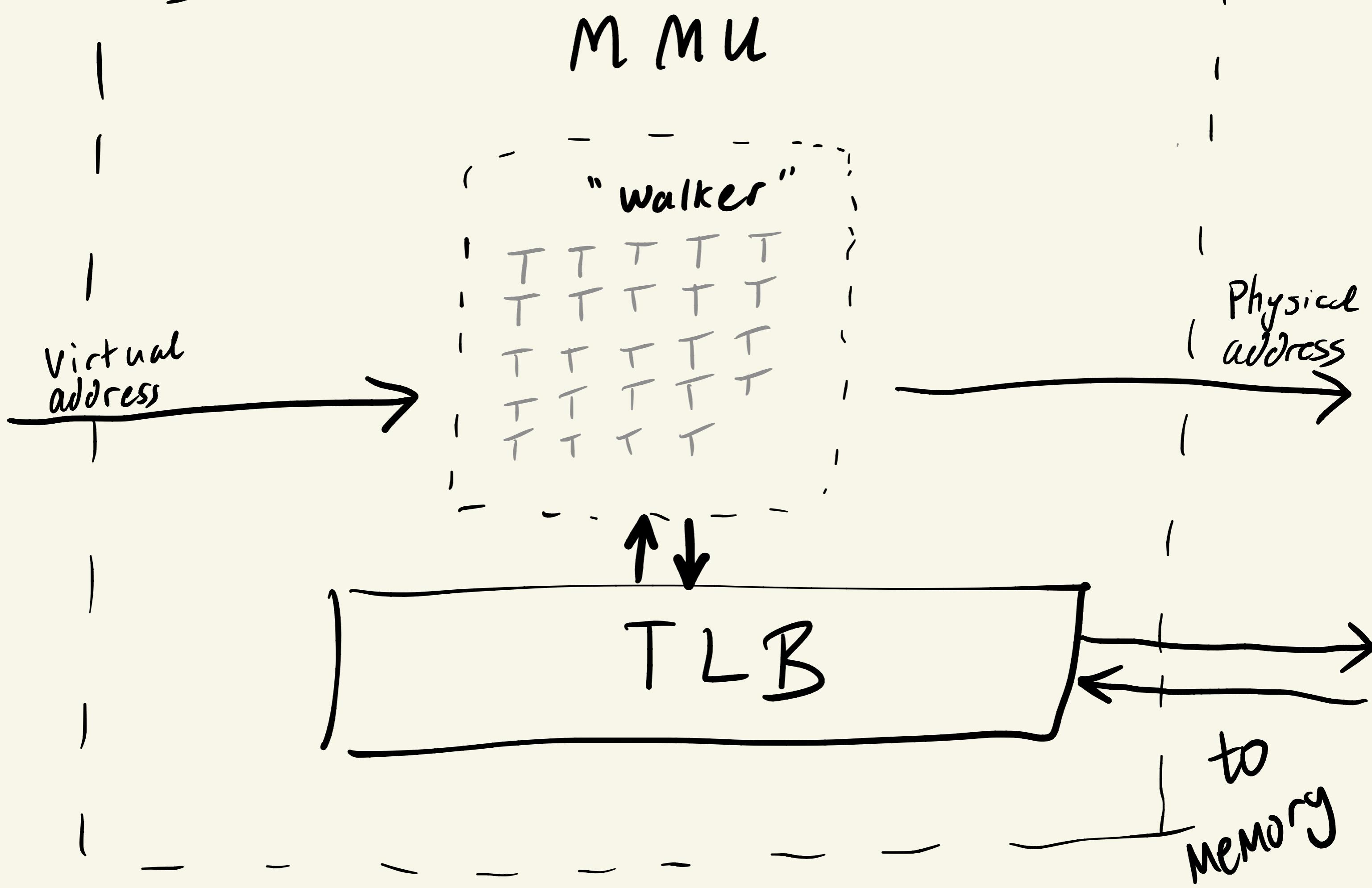


Background

Background







# Issues

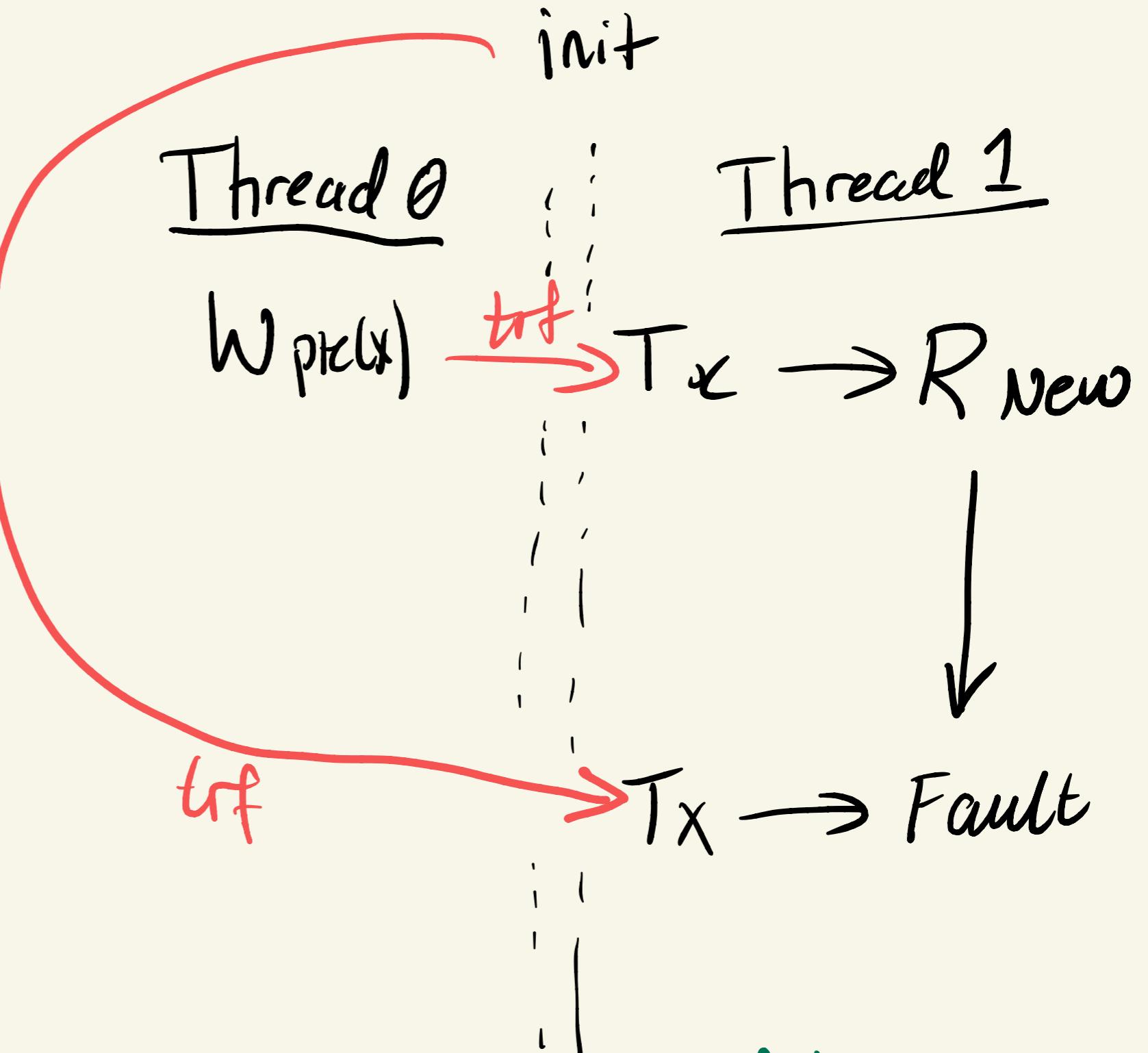
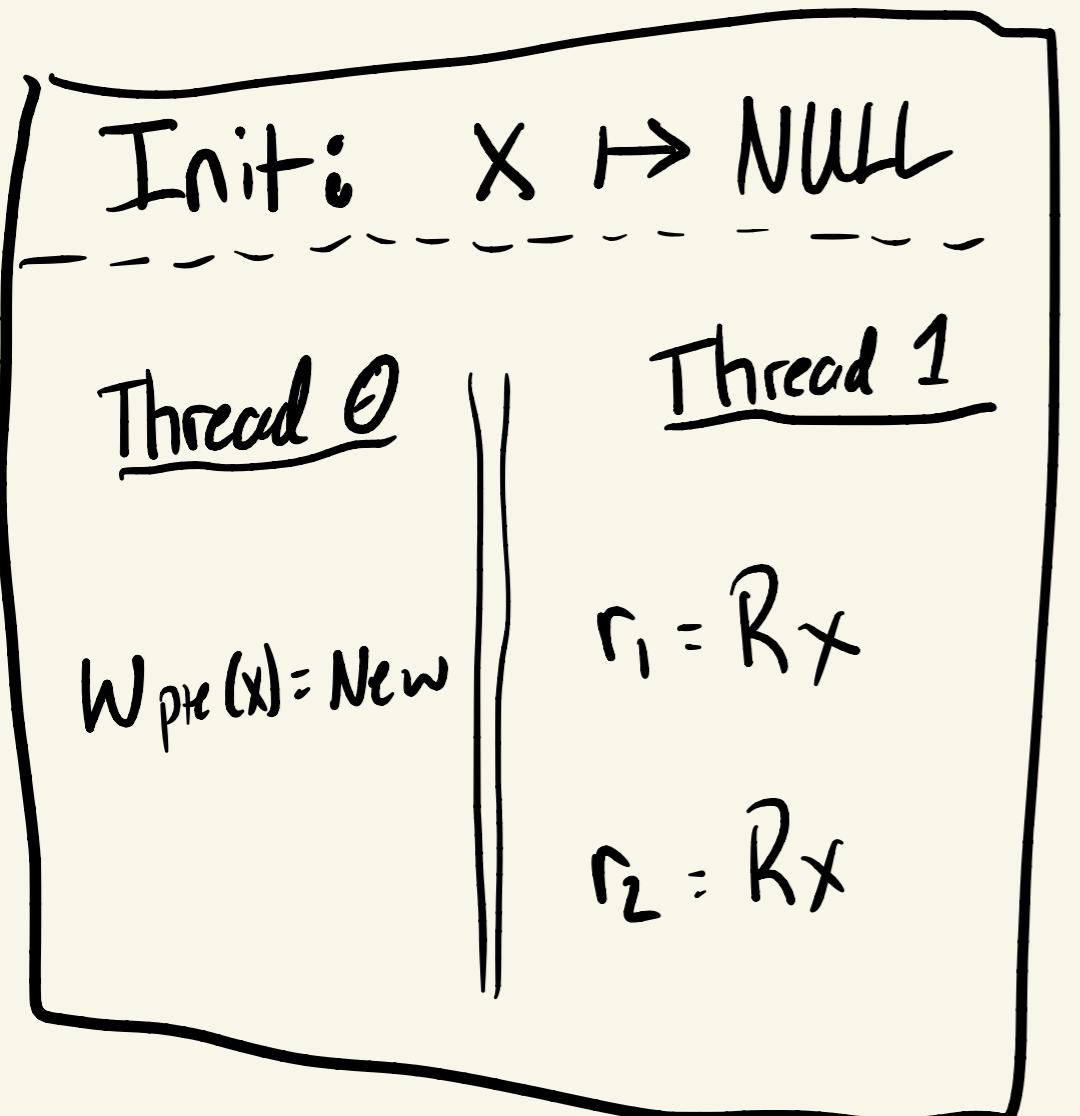
- TLBs
- ARM manual  
11,530 pages

&

# Contributions

- Clarify Spec  
by discussion with Arm
- Extended ARM model
- Protected KVM  
Litmus tests

# Behaviours



Can See 'New'  
then later Stop  
Seeing it?

'Allowed'

# isla – axiomatic

← → C [isla-axiomatic.cl.cam.ac.uk/#](https://isla-axiomatic.cl.cam.ac.uk/#)

Isla W / aarch64-vmsa-strong Litmus file Memory model Sail architecture (aarch64-vmsa) Run test Options Share

W.toml Memory model EventGraph x

```
11     identity 0x1000 with code;
12 """
13
14 [thread.0]
15 init = {}
16 code = """
17     STR X0,[X1]
18 """
19
20 [thread.0.reset]
21 R0 = "desc3(y, page_table_base)"
22 R1 = "pte3(x, page_table_base)"
23
24 [thread.1]
25 init = {}
26 code = """
27     LDR X2,[X1]
28     MOV X0,X2
29     LDR X2,[X3]
30 """
31
32 [thread.1.reset]
33 R1 = "x"
34 R3 = "x"
35 VBAR_EL1 = "extz(0x1000, 64)"
36
37 "PSTATE.SP" = "0b0"
38 "PSTATE.EL" = "0b00"
39
40 [section.thread1_el1_handler]
41 address = "0x1400"
42 code = """
43     MOV X2,#0
44
45     MRS X13,ELR_EL1
46     ADD X13,X13,#4
47     MSR ELR_EL1,X13
48     ERET
49 """
50
51 [final]
52 assertion = "1:X0=1 & 1:X2=0"
```

EventGraph x Relations 1 of 1

The event graph illustrates a race condition between two threads. Thread 0 performs a store operation (a: str x0, [x1]: W s1:l3pte(x) = s1:l3desc(y)) on memory location x1. Thread 1 then attempts to load from the same memory location (a1: T s2:l3pte(x)). This leads to a fault (b2: ldr x2, [x3]: Fault). The initial state is shown as an oval. Transitions are labeled 'trf' (Thread 0 to Thread 1) and 'iio' (Inter-Instruction Operations). A final state 'c: eret' is shown at the bottom right.

# Lots of Behaviours

(Not mentioned here)

- Physical/virtual Coherence
  - Micro TLBs
- Spontaneous translation table walks
  - Multi-copy atomicity
  - ETS
- Coherent TLB fills
  - Break-before-make
- Re-ordered translations
  - write forwarding
- In-order walks
  - TLB Shootdown
- Speculative translations
  - ASIDs and VMIDs

# Axiomatic Model

```

let tlb-affects =
  ...
let TLB_barrier =
  ([TLBI] ; tlb-affects ; [T] ; tfr ; [W])^-1 & wco
let maybe_TLB_cached =
  ([T] ; trf^-1 ; wco ; [TLBI-S1]) & tlb-affects^-1
let tcachel = [T & Stage1] ; tfr ; TLB_barrier
let tcache2 = [T & Stage2] ; tfr ; TLB_barrier

let speculative =
  ctrl
  | addr; po
  | [T] ; instruction-order
(* translation-ordered-before *)
let tob =
  [T_f] ; tfre
  | ([T_f] ; tfri) & (po ; [DSB.SY] ; instruction-order)^-1
  | [T] ; iio ; [R|W] ; po ; [W]
  | speculative ; trfi
(* observed by *)
let obs = rfe | fr | wco
  | trfe
(* ordered-before TLBI and translate *)
let obtlbi_translate =
  tcachel
  | tcache2 & (iio^-1 ; [T & Stage1] ; trf^-1 ; wco^-1)
  | (tcache2 ; wco? ; [TLBI-S1]) & (iio^-1 ; [T & Stage1] ;
    maybe_TLB_cached)

(* ordered-before TLBI *)
let obtlbi =
  obtlbi_translate
  | [R|W|Fault] ; iio^-1 ; (obtlbi_translate & ext) ; [TLBI]

(* context-change ordered-before *)
let ctxob =
  speculative ; [MSR]
  | [CSE] ; instruction-order
  | [ContextChange] ; po ; [CSE]
  | speculative ; [CSE]
  | po ; [ERET] ; instruction-order ; [T]

(* ordered-before a translation fault *)
let obfault =

```

```

data ; [Fault & IsFromW]
| speculative ; [Fault & IsFromW]
| [dmbst] ; po ; [Fault & IsFromW]
| [dmbld] ; po ; [Fault & (IsFromW|IsFromR)]
| [A|Q] ; po ; [Fault & (IsFromW | IsFromR)]
| [R|W] ; po ; [Fault & IsFromW & IsReleaseW]
(* ETS-ordered-before *)
let obETS =
  (obfault ; [Fault]) ; iio^-1 ; [T_f]
  | ([TLBI] ; po ; [dsb] ; instruction-order ; [T]) & tlb-affects

(* dependency-ordered-before *)
let dob =
  addr | data
  | speculative ; [W]
  | addr; po; [W]
  | (addr | data); rfi
  | (addr | data); trfi
(* atomic-ordered-before *)
let aob = rmw
  | [range(rmw)]; rfi; [A | Q]

(* barrier-ordered-before *)
let bob = [R] ; po ; [dmbld]
  | [W] ; po ; [dmbst]
  | [dmbst]; po; [W]
  | [dmbld]; po; [R|W]
  | [L]; po; [A]
  | [A | Q]; po; [R | W]
  | [R | W]; po; [L]
  | [F | C]; po; [dsbsy]
  | [dsb] ; po

(* Ordered-before *)
let ob = (obs | dob | aob | bob | iio | tob | obtlbi | ctxob |
  , - obfault | obETS)^+
  | (* Internal visibility requirement *)
  | acyclic po-loc | fr | co | rf as internal
  | (* External visibility requirement *)
  | irreflexive ob as external
  | (* Atomic requirement *)
  | empty rmw & (fre; coe) as atomic
  | (* Writes cannot forward to po-future translates *)
  | acyclic (po-pa | trfi) as translation-internal
  )

```

Relations

Axioms

# PKVM Examples

pkvm.unshare-page (broken):

// unmap memory

pte(i) := NULL

// continue guest OS

return

OS:

:

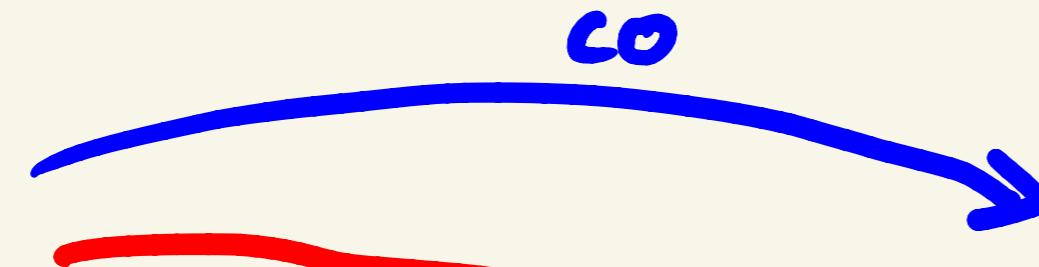
:

:

r<sub>0</sub> := [v]

Must trap?

$w \text{ pte}(i)$



$W \quad \text{pte}(i) = 0$

A memory map diagram showing a 5x5 grid of T's. A red arrow labeled "brk" points to the bottom-right cell, which contains a circled 'T'. A green oval highlights this circled 'T'.

T	T	T	T	T
T	T	T	T	T
T	T	T	T	T
T	T	T	T	T
T	T	T	+	T

{ what is sufficient here  
to force the read to fault?

return

R pa1

PKVM :

$\text{pte}(i) := \text{NULL}$

DSB

TLBI IPA, i

DSB

TLBI S1

DSB

Return

OS1 :

⋮

$r_0 := [v]$

Now  
must  
trap!

```

let tlb-affects =
  ...

let TLB_barrier =
  ([TLBI] ; tlb-affects ; [T] ; tfr ; [W])^-1
  & wco

let maybe_TLB_cached =
  ([T] ; trf^-1 ; wco ; [TLBI-S1]) & tlb-affects^-1

let tcache1 = [T & Stage1] ; tfr ; TLB_barrier
let tcache2 = [T & Stage2] ; tfr ; TLB_barrier

(* ordered-before TLBI and translate *)
let obtlbi_translate =
  ...
  | (tcache2 ; wco? ; [TLBI-S1])
    & (iio^-1 ; [T & Stage1] ; maybe_TLB_cached)

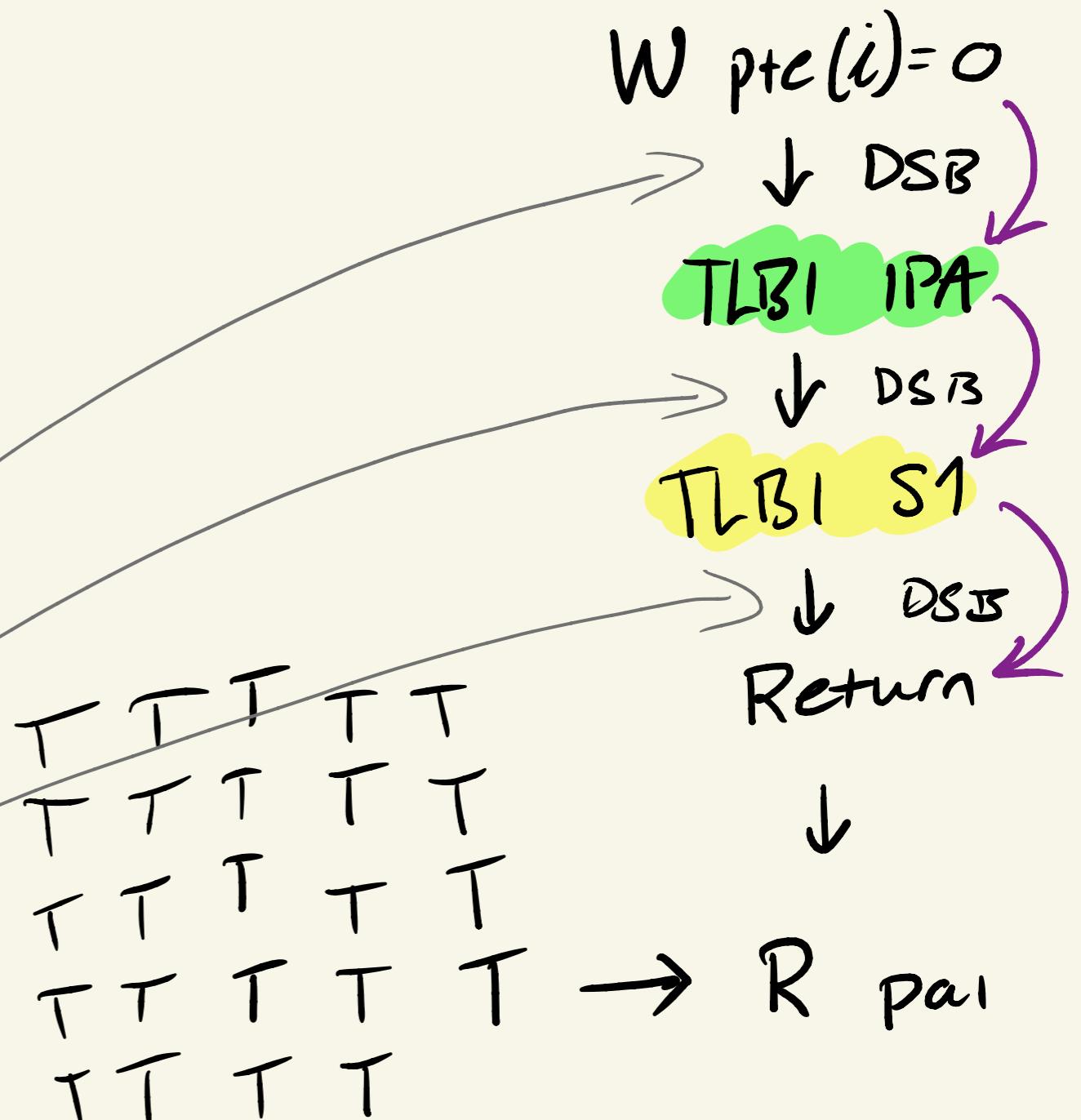
(* ordered-before TLBI *)
let obtlbi =
  obtlbi_translate
  | ...

(* context-change ordered-before *)
let ctxob =
  ...
  | [CSE] ; instruction-order
  | po ; [ERET] ; instruction-order ; [T]

(* barrier-ordered-before *)
let bob =
  ...
  | [F | C] ; po ; [dsbsy]
  | [dsb] ; po

(* Ordered-before *)
let ob = (... | bob | iio | obtlbi | ctxob)^+
irreflexive ob as external

```



```

let tlb_might_affect =
  [ TLBI-S1 & ~TLBI-S2 & TLBI-VA & TLBI-ASID & TLBI-VMID] ; (same-va-page & same-asid & same-vmid) ; [T & Stage1]
| [ TLBI-S1 & ~TLBI-S2 & ~TLBI-VA & TLBI-ASID & TLBI-VMID] ; (same-asid & same-vmid) ; [T & Stage1]
| [ TLBI-S1 & ~TLBI-S2 & ~TLBI-VA & ~TLBI-ASID & TLBI-VMID] ; same-vmid ; [T & Stage1]
| [~TLBI-S1 & TLBI-S2 & TLBI-IPA & ~TLBI-ASID & TLBI-VMID] ; (same-ipa-page & same-vmid) ; [T & Stage2]
| [~TLBI-S1 & TLBI-S2 & ~TLBI-IPA & ~TLBI-ASID & TLBI-VMID] ; same-vmid ; [T & Stage2]
| [ TLBI-S1 & TLBI-S2 & ~TLBI-IPA & ~TLBI-ASID & TLBI-VMID] ; same-vmid ; [T]
| ( TLBI-S1 & ~TLBI-IPA & ~TLBI-ASID & ~TLBI-VMID) * (T & Stage1)
| ( TLBI-S2 & ~TLBI-IPA & ~TLBI-ASID & ~TLBI-VMID) * (T & Stage2)

let tlb-affects =
  [TLBI-IS] ; tlb_might_affect
| ([~TLBI-IS] ; tlb_might_affect) & int

```

T T T T  
 T T T T  
 T T T T  
 T T T T  
 T T T T  
 T T T T

T T T  
 T T T  
 T T T

T T T  
 T T T  
 T T T

T T T

W  $p_{tc}(i) = 0$

TLBI IPA

TLBI S1

return

R pa1

DSB

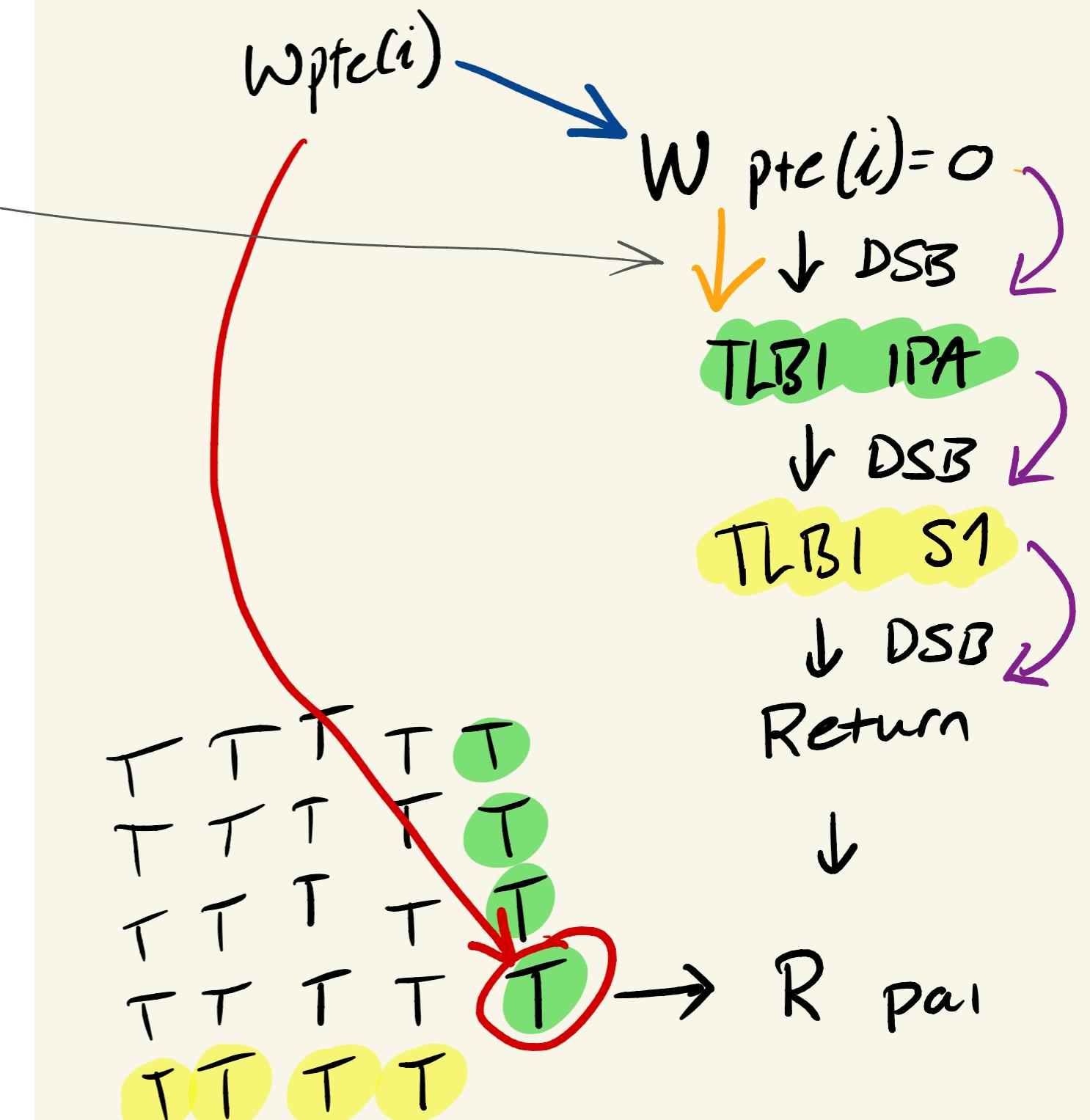
DSB

DSB

```

let tlb-affects =
  ...
let TLB_barrier =
  ([TLBI] ; tlb-affects ; [T] ; tfr ; [W])^-1
  & wco
let maybe_TLB_cached =
  ([T] ; trf^-1 ; wco ; [TLBI-S1]) & tlb-affects^-1
let tcache1 = [T & Stage1] ; tfr ; TLB_barrier
let tcache2 = [T & Stage2] ; tfr ; TLB_barrier
(* ordered-before TLBI and translate *)
let obtlbi_translate =
  ...
  | (tcache2 ; wco? ; [TLBI-S1])
    & (iio^-1 ; [T & Stage1] ; maybe_TLB_cached)
(* ordered-before TLBI *)
let obtlbi =
  obtlbi_translate
  | ...
(* context-change ordered-before *)
let ctxob =
  ...
  | [CSE] ; instruction-order
  | po ; [ERET] ; instruction-order ; [T]
(* barrier-ordered-before *)
let bob =
  ...
  | [F | C] ; po ; [dsbsy]
  | [dsb] ; po
(* Ordered-before *)
let ob = (... | bob | iio | obtlbi | ctxob)^+
irreflexive ob as external

```



```

let tlb-affects =
  ...
let TLB_barrier =
  ([TLBI] ; tlb-affects ; [T] ; tfr ; [W])^-1
  & wco

let maybe_TLB_cached =
  ([T] ; trf^-1 ; wco ; [TLBI-S1]) & tlb-affects^-1

let tcache1 = [T & Stage1] ; tfr ; TLB_barrier
let tcache2 = [T & Stage2] ; tfr ; TLB_barrier

(* ordered-before TLBI and translate *)
let obtlbi_translate =
  ...
  | (tcache2 ; wco? ; [TLBI-S1])
    & (iio^-1 ; [T & Stage1] ; maybe_TLB_cached)

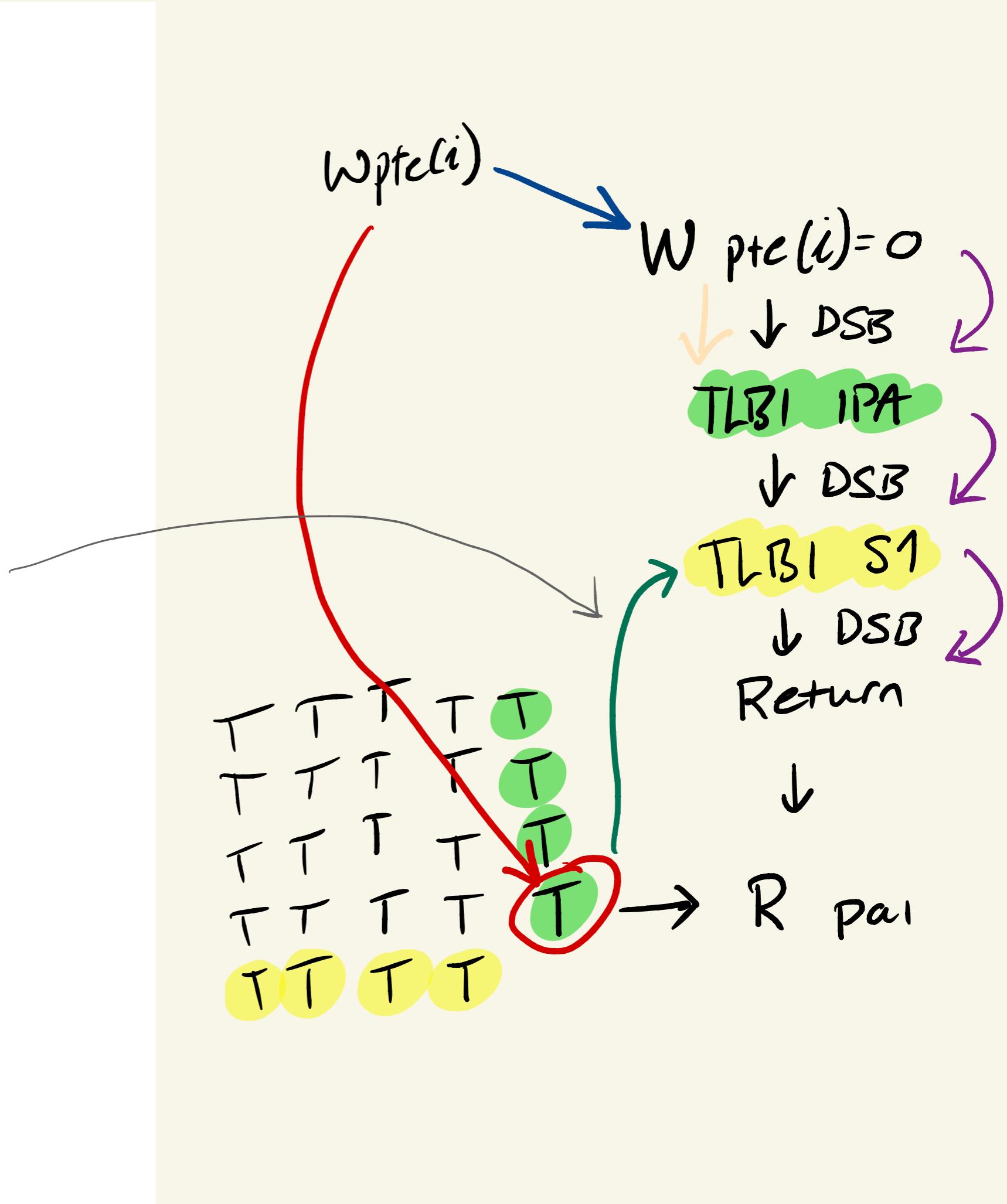
(* ordered-before TLBI *)
let obtlbi =
  obtlbi_translate
  | ...

(* context-change ordered-before *)
let ctxob =
  ...
  | [CSE] ; instruction-order
  | po ; [ERET] ; instruction-order ; [T]

(* barrier-ordered-before *)
let bob =
  ...
  | [F | C] ; po ; [dsbsy]
  | [dsb] ; po

(* Ordered-before *)
let ob = (... | bob | iio | obtlbi | ctxob)^+
irreflexive ob as external

```



```

let tlb-affects =
  ...

let TLB_barrier =
  ([TLBI] ; tlb-affects ; [T] ; tfr ; [W])^-1
  & wco

let maybe_TLB_cached =
  ([T] ; trf^-1 ; wco ; [TLBI-S1]) & tlb-affects^-1

let tcache1 = [T & Stage1] ; tfr ; TLB_barrier
let tcache2 = [T & Stage2] ; tfr ; TLB_barrier

(* ordered-before TLBI and translate *)
let obtlbi_translate =
  ...
  | (tcache2 ; wco? ; [TLBI-S1])
    & (iio^-1 ; [T & Stage1] ; maybe_TLB_cached)

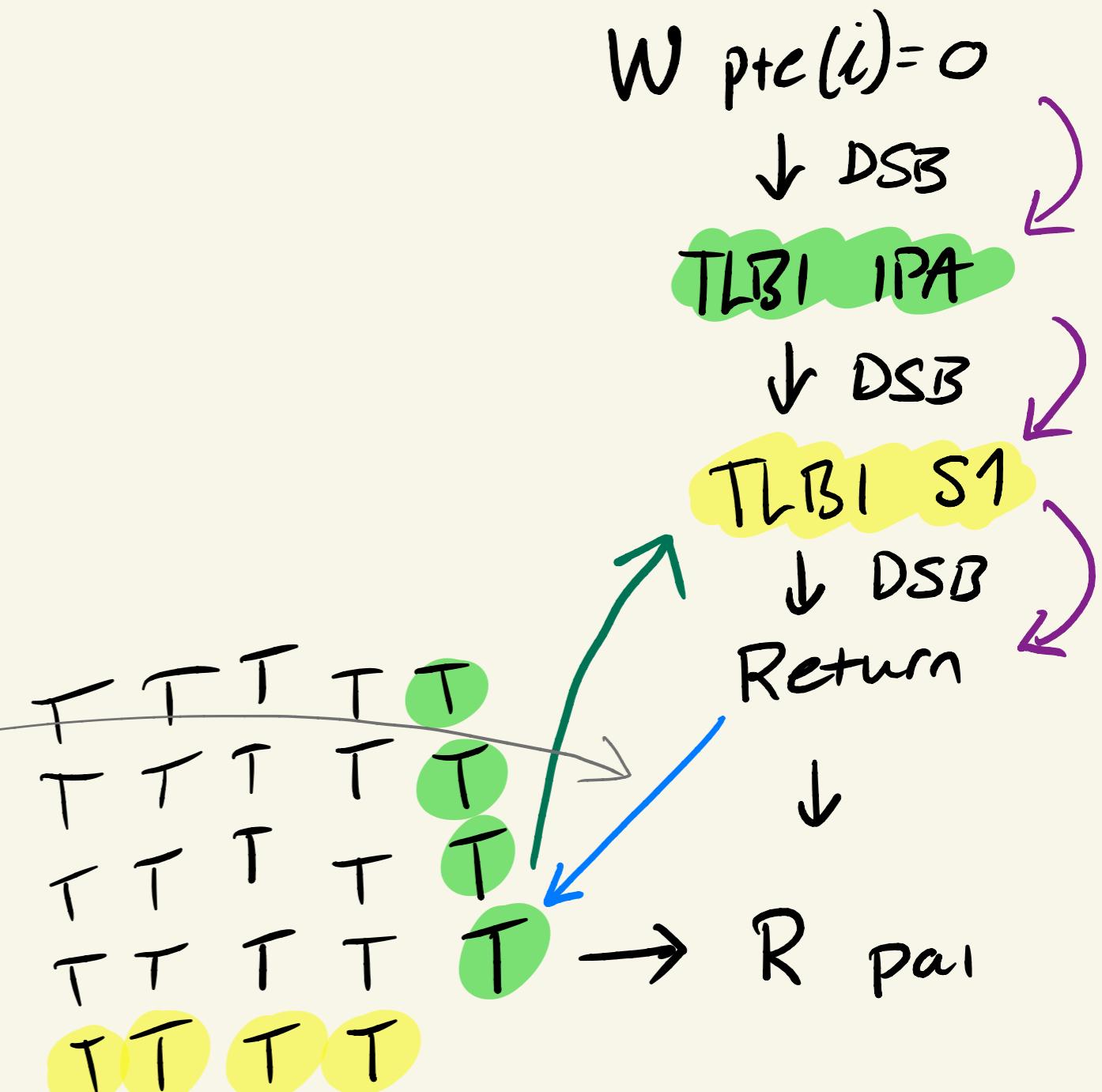
(* ordered-before TLBI *)
let obtlbi =
  obtlbi_translate
  | ...

(* context-change ordered-before *)
let ctxob =
  ...
  | [CSE] ; instruction-order
  | po ; [ERET] ; instruction-order ; [T]

(* barrier-ordered-before *)
let bob =
  ...
  | [F | C] ; po ; [dsbsy]
  | [dsb] ; po

(* Ordered-before *)
let ob = (... | bob | iio | obtlbi | ctxob)^+
irreflexive ob as external

```



```

let tlb-affects =
  ...

let TLB_barrier =
  ([TLBI] ; tlb-affects ; [T] ; tfr ; [W])^-1
  & wco

let maybe_TLB_cached =
  ([T] ; trf^-1 ; wco ; [TLBI-S1]) & tlb-affects^-1

let tcache1 = [T & Stage1] ; tfr ; TLB_barrier
let tcache2 = [T & Stage2] ; tfr ; TLB_barrier

(* ordered-before TLBI and translate *)
let obtlbi_translate =
  ...
  | (tcache2 ; wco? ; [TLBI-S1])
    & (iio^-1 ; [T & Stage1] ; maybe_TLB_cached)

(* ordered-before TLBI *)
let obtlbi =
  obtlbi_translate
  | ...

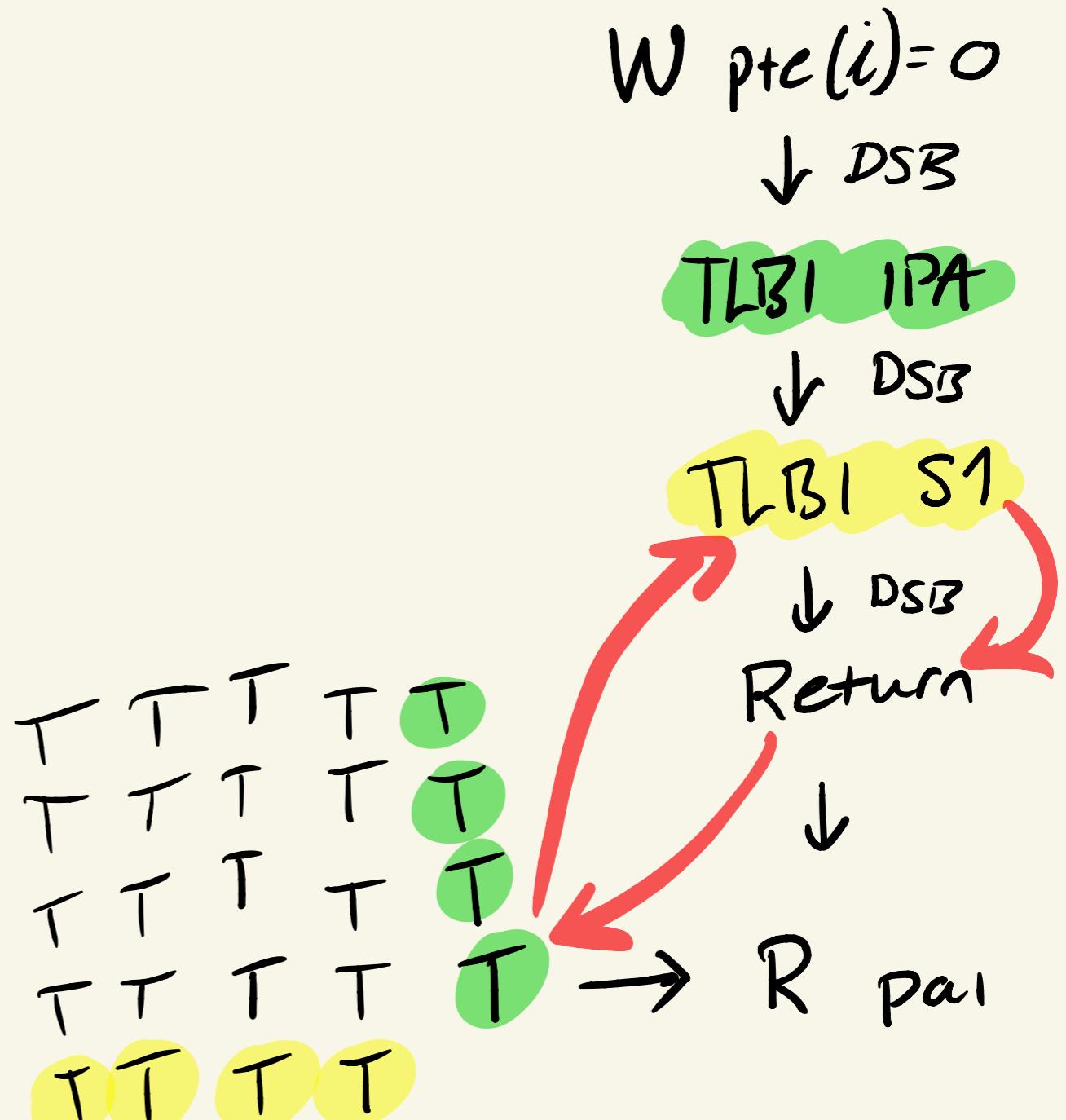
(* context-change ordered-before *)
let ctxob =
  ...
  | [CSE] ; instruction-order
  | po ; [ERET] ; instruction-order ; [T]

(* barrier-ordered-before *)
let bob =
  ...
  | [F | C] ; po ; [dsbsy]
  | [dsb] ; po

(* Ordered-before *)
let ob = (... | bob | iio | obtlbi | ctxob)^+

```

irreflexive ob as external



# Recap

- ① Clarified Architecture
- ② Axiomatize model  
+ tooling
- ③ pkVM tests

Future

use the model !

More Arm features

... more architectures?