

A case against semantic dependencies

***Ori Lahav*, Tel Aviv University**

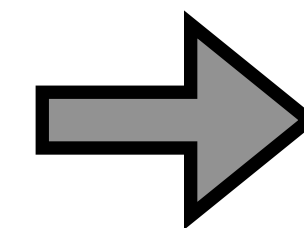
Based on joint discussions with *Minki Cho, Chung-Kil Hur, Sung-Hwan Lee, Ben Simner*

Imagine a new model for C/C++ that works in three steps:

- Step 1: “calculate” a set S of **candidate program execution graphs**
- Step 2: given S , derive **semantic dependency** (**sdep**) for each graph
- Step 3: apply the **consistency predicate** from the C/C++ standard

Definition 1. An execution G is called *RC11-consistent* if it is complete and the following hold:

- $hb; eco^?$ is irreflexive. (COHERENCE)
- $rmw \cap (rb; mo) = \emptyset$. (ATOMICITY)
- psc is acyclic. (SC)
- $po \cup rf$ is acyclic. (NO-THIN-AIR)



sdep \cup **rf** is acyclic.



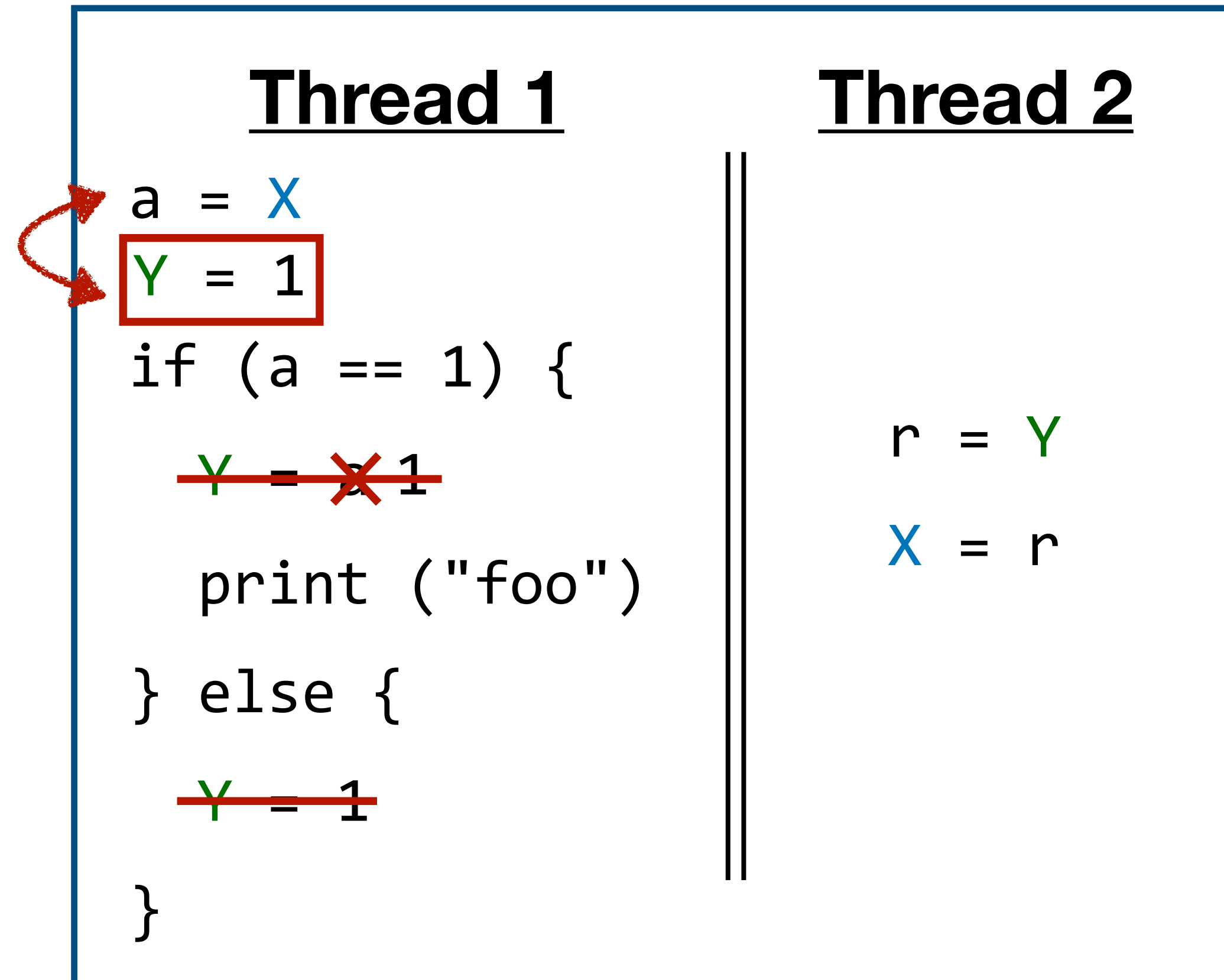
Imagine a new model for C/C++ that works in three steps:

- Step 1: “calculate” a set S of **candidate program execution graphs**
- Step 2: given S, derive **semantic dependency** (**sdep**) for each graph
- Step 3: apply the **consistency predicate** from the C/C++ standard

I believe this approach can't work. I argue via example that:

- *Step 2 cannot be thread-local*
- *Step 2 has to be aware of the consistency predicate in step 3*

Can “foo” be printed?



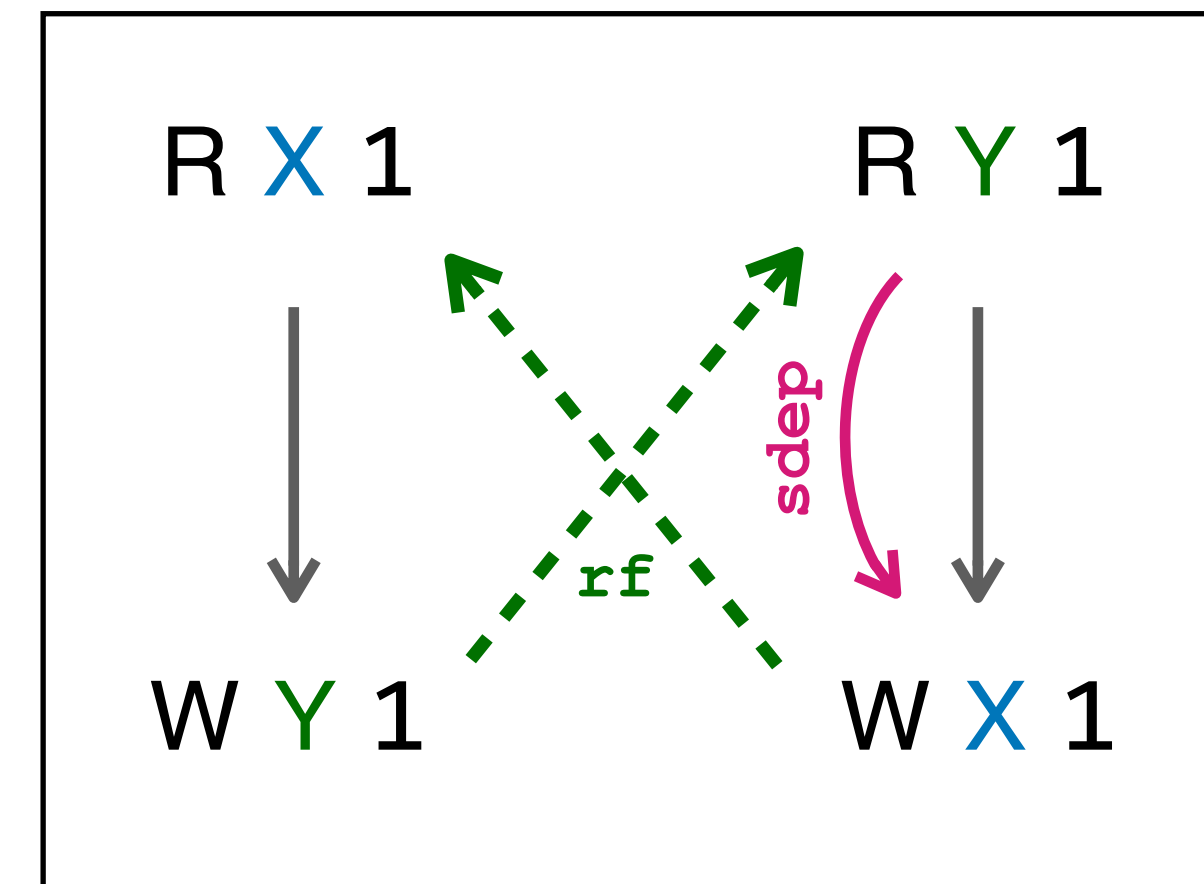
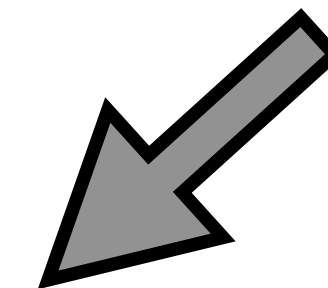
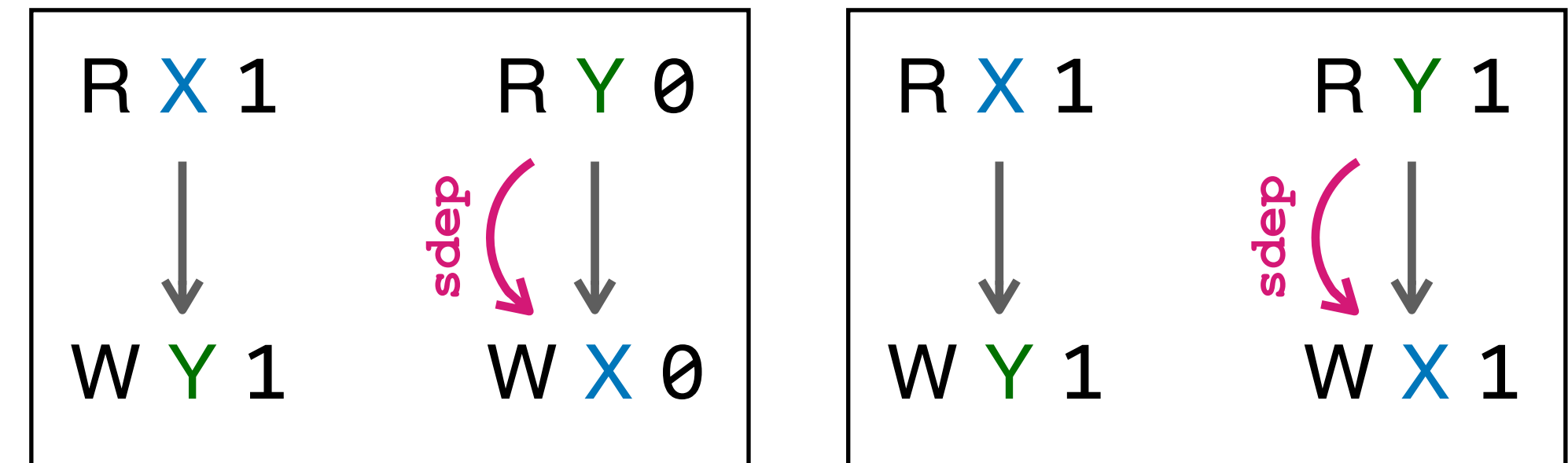
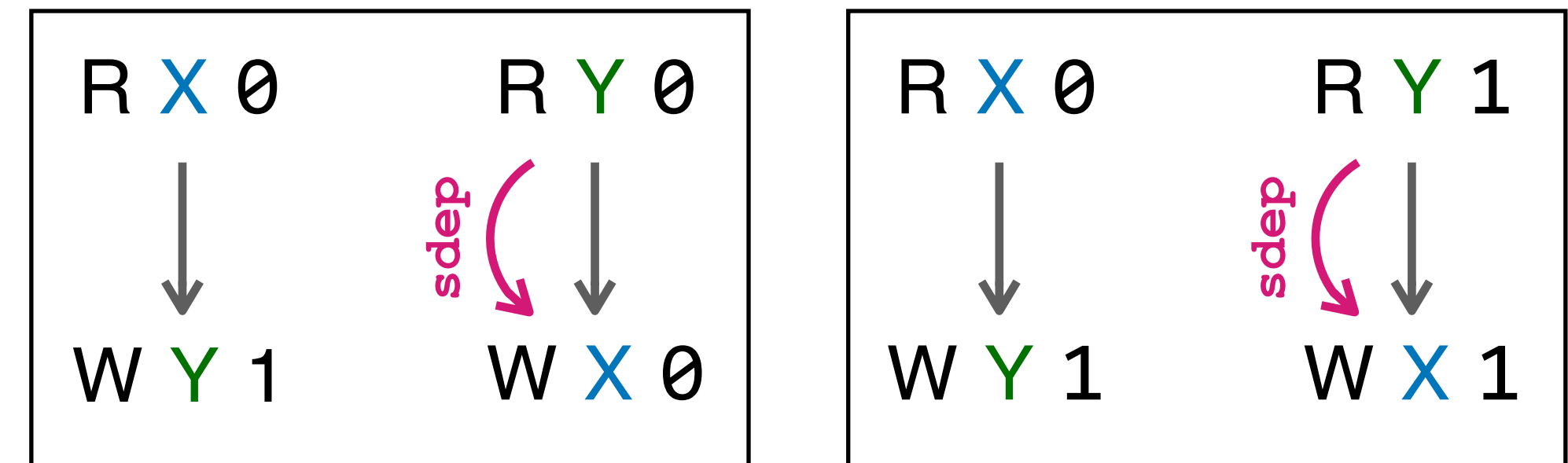
sdep approach

Thread 1

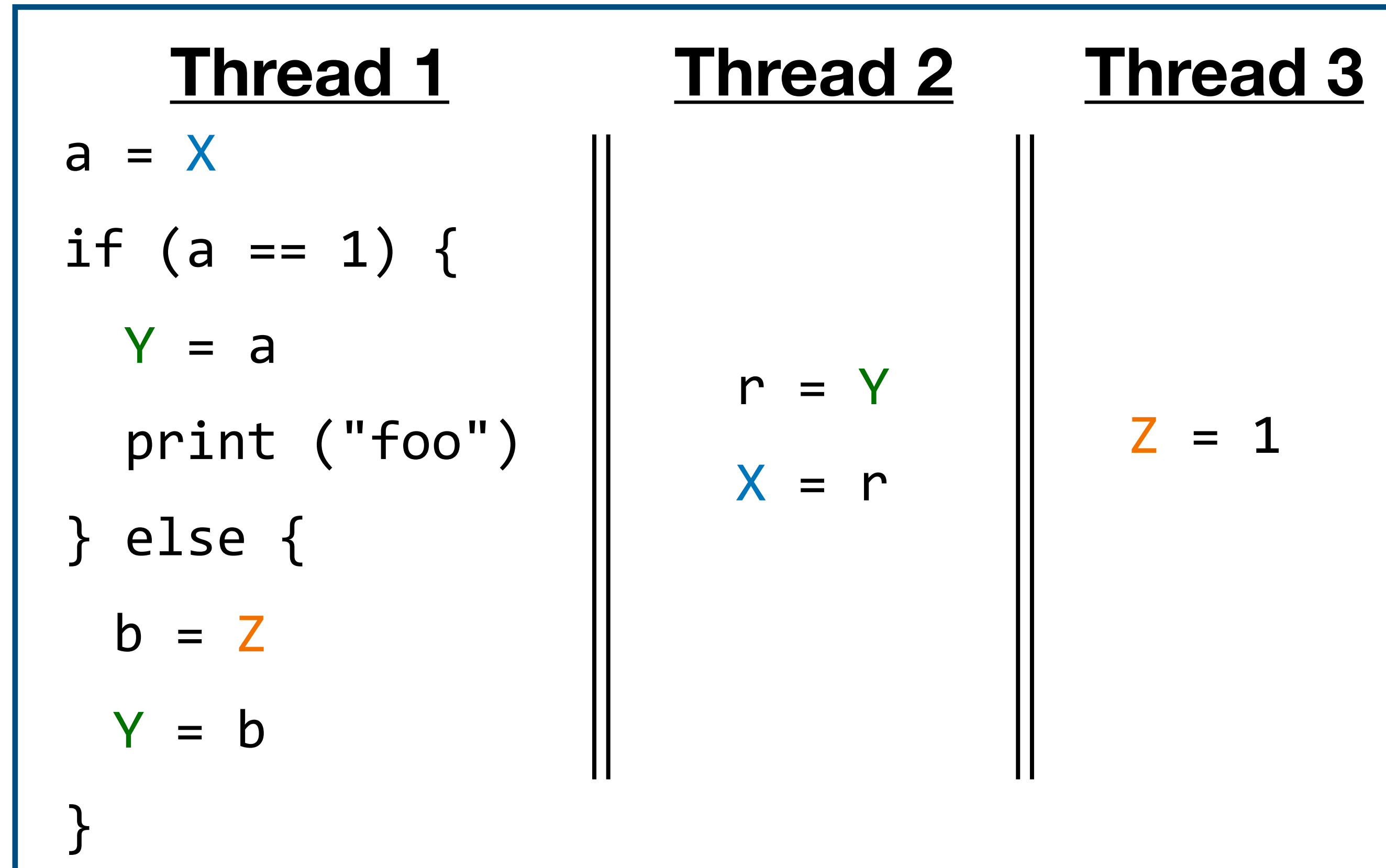
```
a = X
if (a == 1) {
  Y = a
  print ("foo")
} else {
  Y = 1
}
```

Thread 2

```
r = Y
X = r
```



Main example



Thread 1

```

a = X
if (a == 1) {
  Y = a
  print ("foo")
} else {
  b = Z
  Y = b
}

```

Thread 2

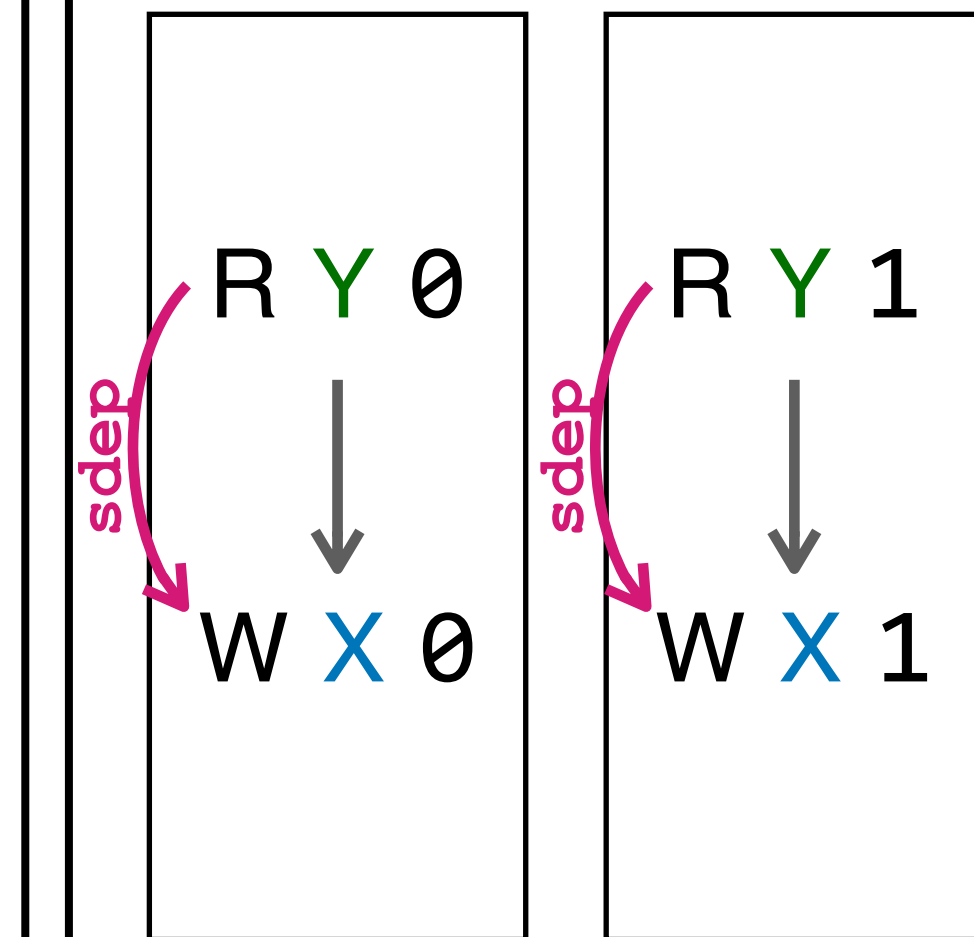
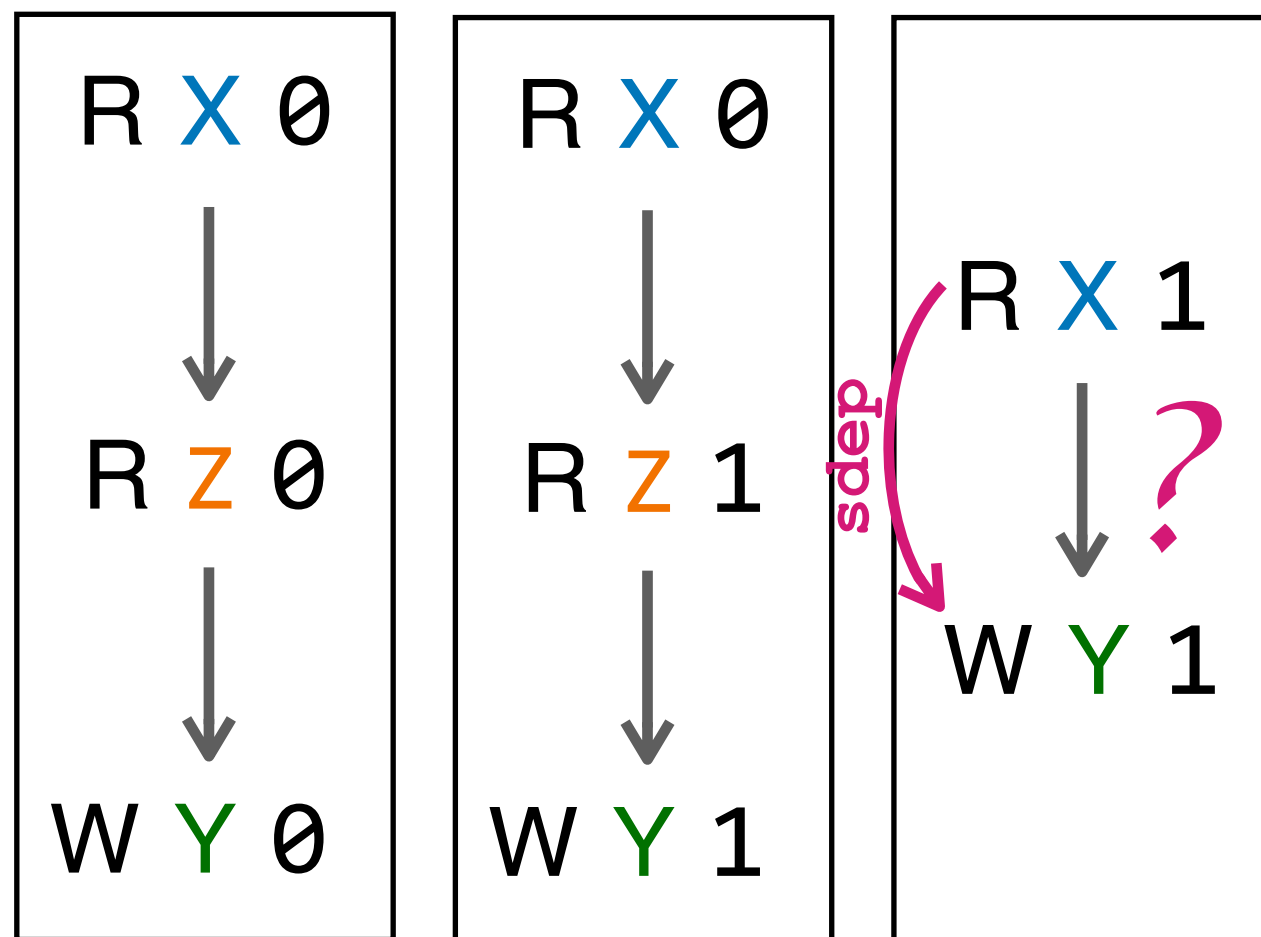
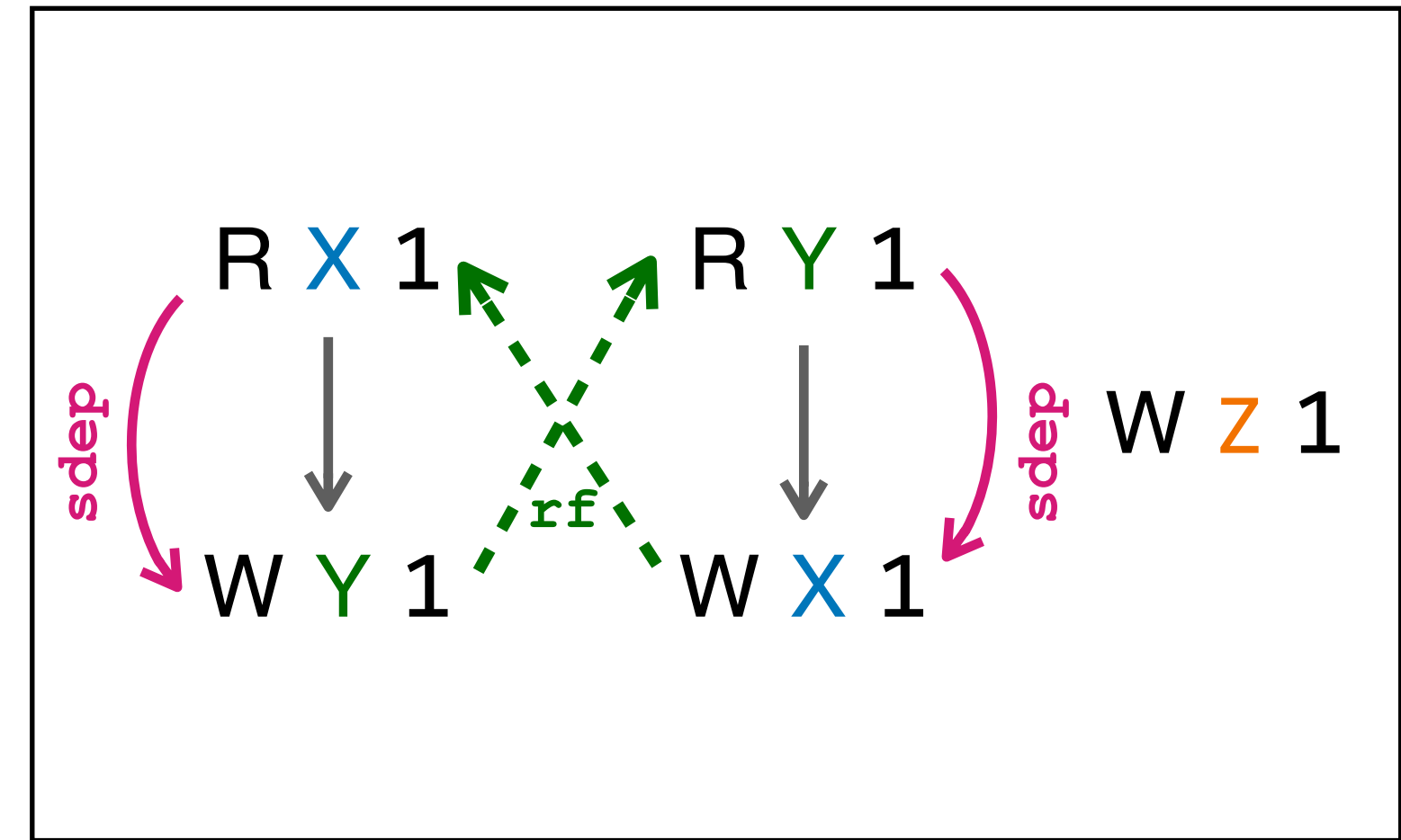
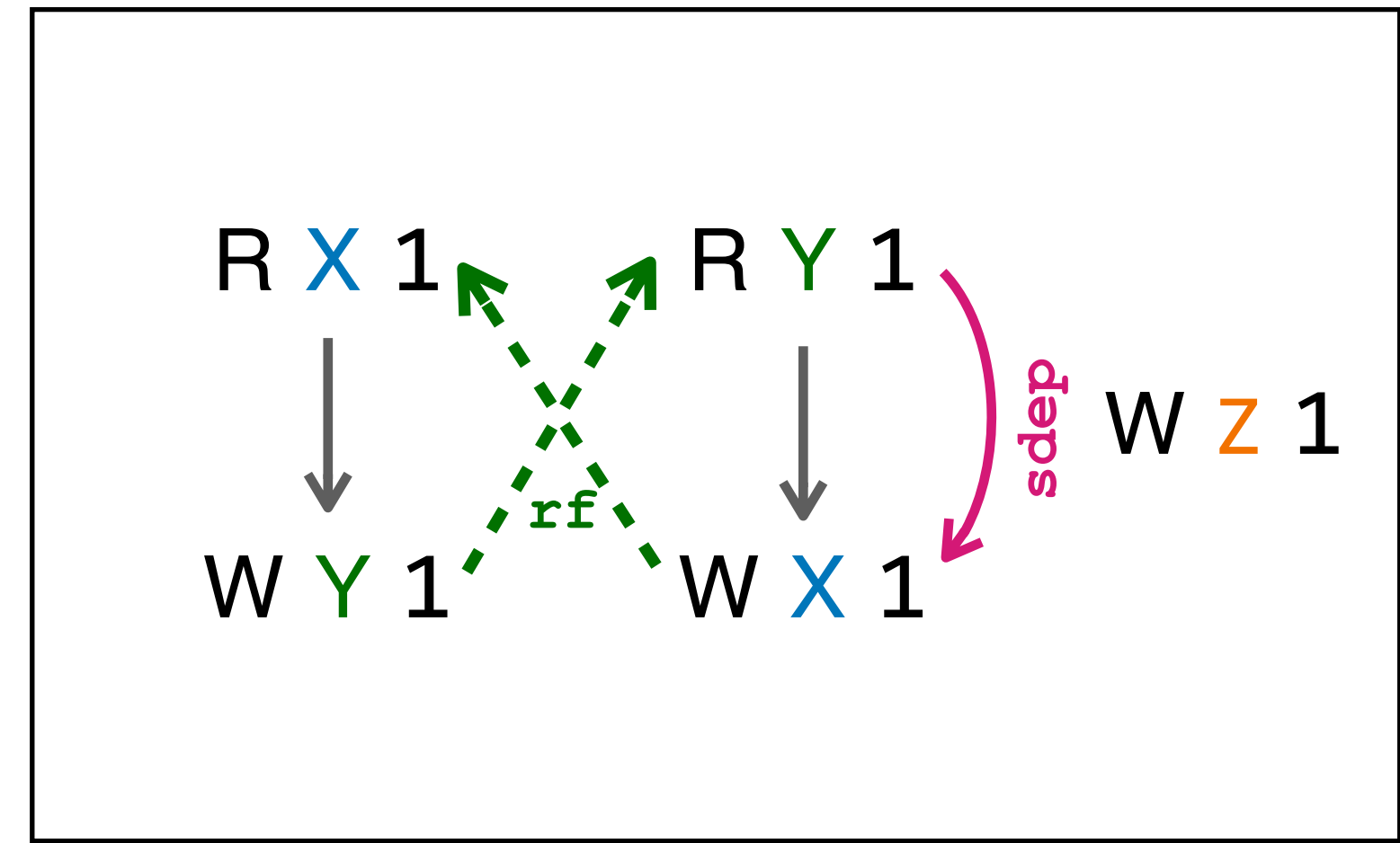
```

r = Y
X = r

```

Thread 3

```
Z = 1
```



Thread 1

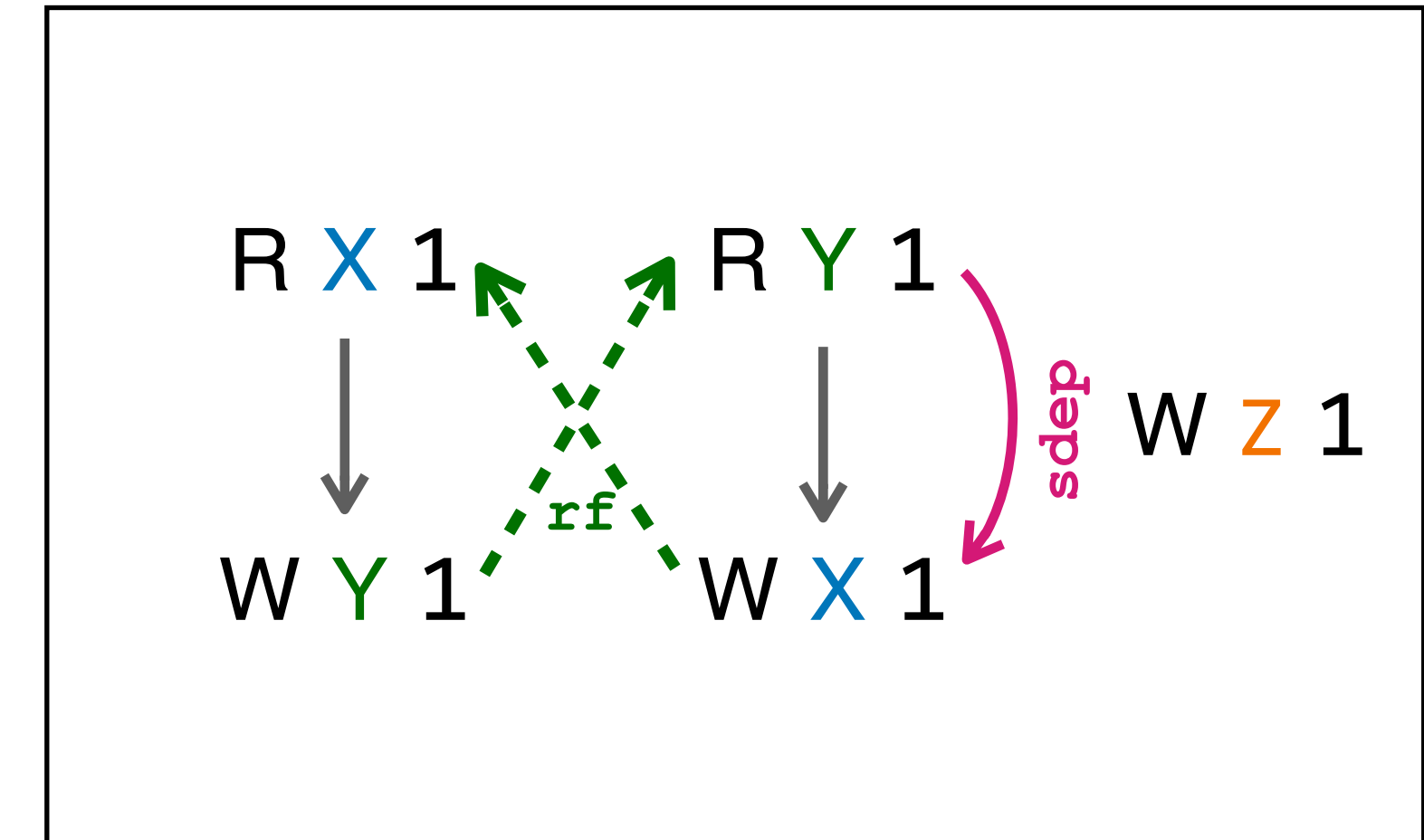
```
a = X
if (a == 1) {
    Y = a
    print ("foo")
} else {
    b = Z
    Y = b
}
```

Thread 2

```
r = Y
X = r
```

Thread 3

```
Z = 1
```



- Printing “foo” **has to be allowed**, assuming we allow compilers to:

- Introduce redundant loads
- Forward load across atomics:

*Both are performed by LLVM/GCC on non-atomics
(Z can be easily made non-atomic)*

c = Z; a = X; b = Z → c = Z; a = X; b = c

Thread 1

```
a = X
if (a == 1) {
  Y = a
  print ("foo")
} else {
  b = Z
  Y = b
}
```

Thread 1

```
c = Z
if (c == 1) {
  a = X
  if (a == 1) {
    Y = a
    print ("foo")
  } else {
    b = Z
    Y = b
  }
} else {
  ...
}
```

Thread 1

```
c = Z
if (c == 1) {
  a = X
  if (a == 1) {
    Y = a
    print ("foo")
  } else {
    b = c
    Y = b
  }
} else {
  ...
}
```

Thread 1

```
c = Z
if (c == 1) {
  a = X
  if (a == 1) {
    Y = 1
    print ("foo")
  } else {
    b = 1
    Y = 1
  }
} else {
  ...
}
```

• • •

Thread 1

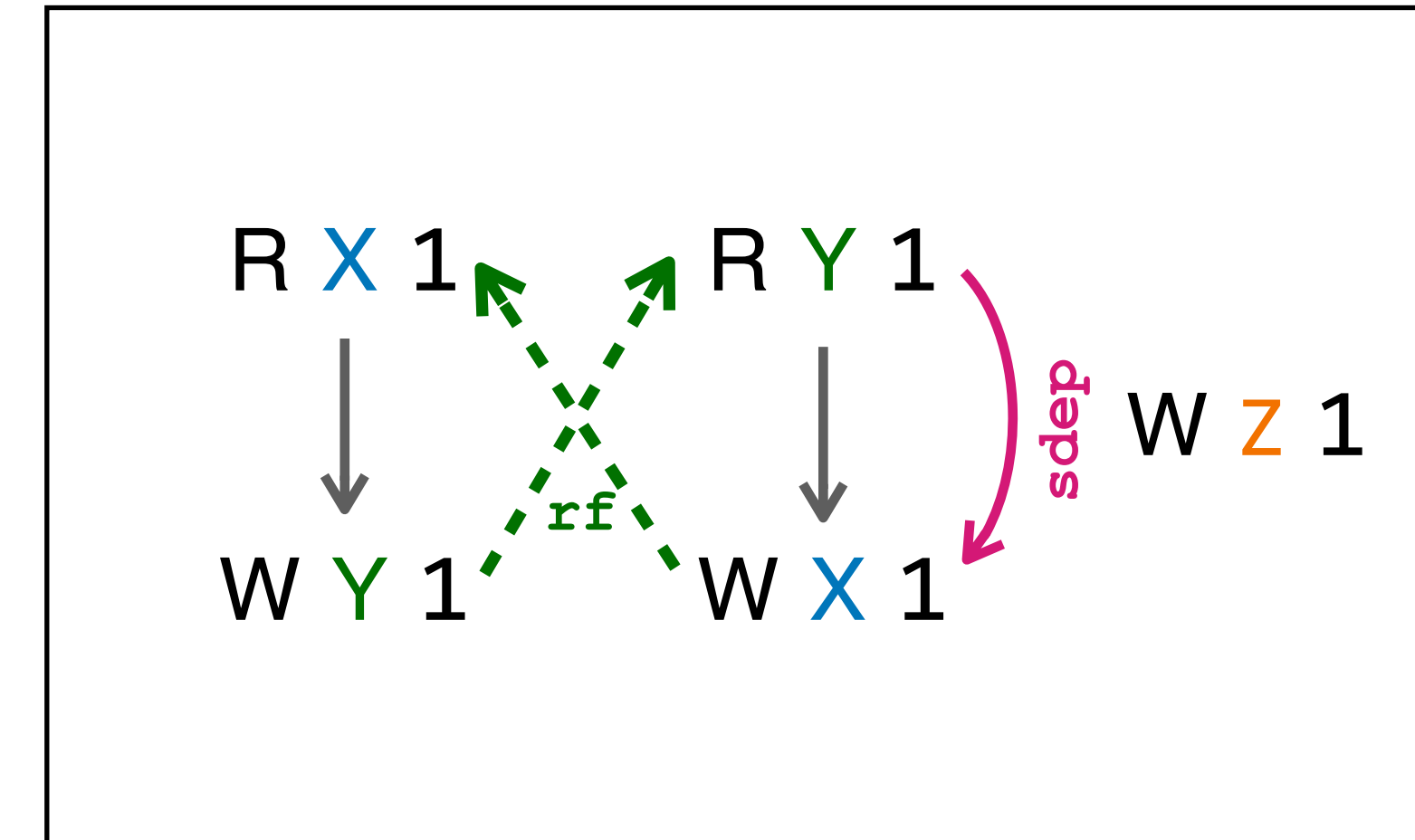
```
a = X
if (a == 1) {
  Y = a
  print ("foo")
} else {
  b = Z
  Y = b
}
```

Thread 2

```
r = Y
X = r
```

Thread 3

```
Z = 1
```



- Printing “foo” **has to be allowed**, assuming we allow compilers to:

- Introduce redundant loads
- Forward load across atomics:

*Both are performed by LLVM/GCC on non-atomics
(Z can be easily made non-atomic)*

c = Z; a = X; b = Z → c = Z; a = X; b = c

Thread 1

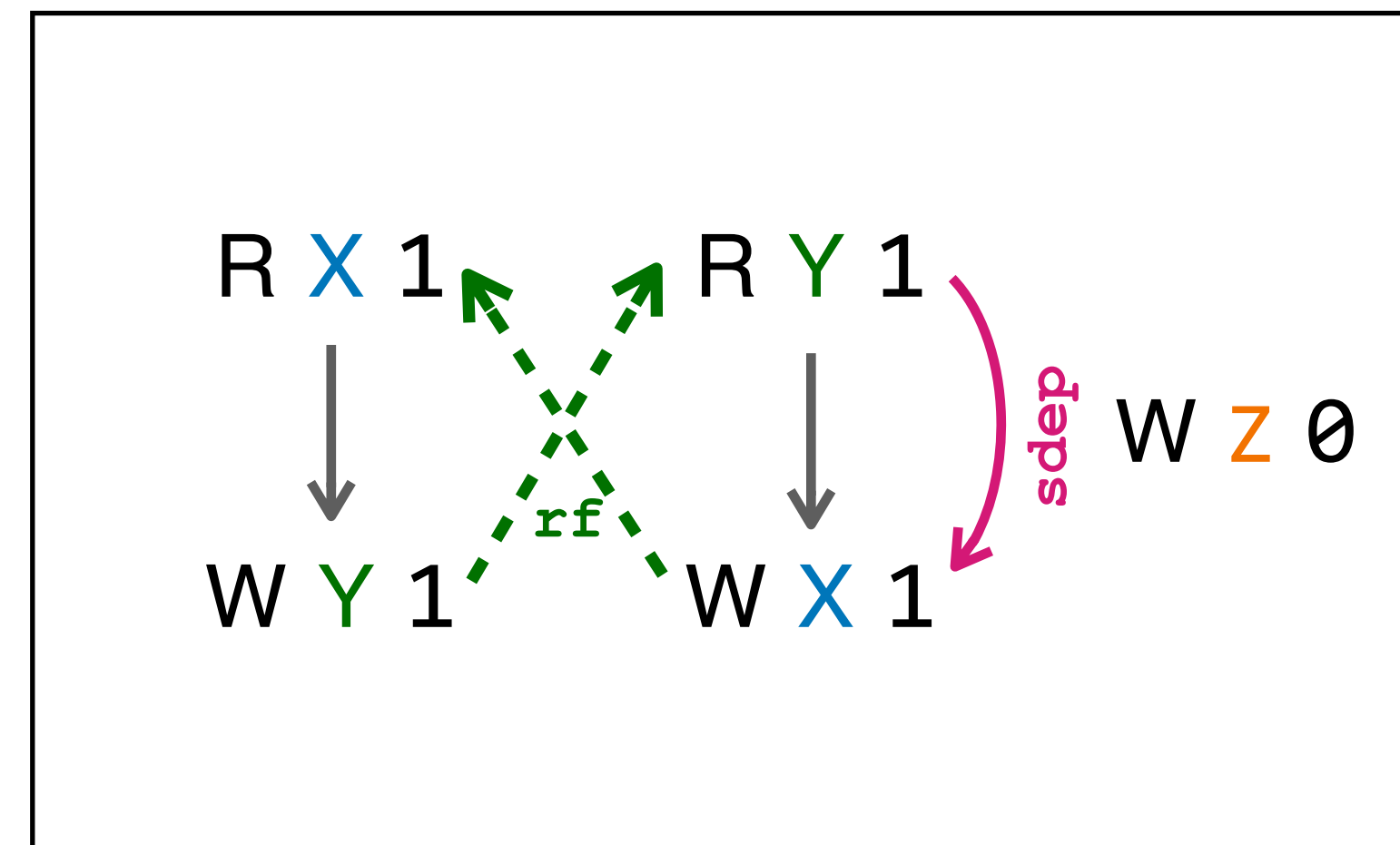
```
a = X
if (a == 1) {
  Y = a
  print ("foo")
} else {
  b = Z
  Y = b
}
```

Thread 2

```
r = Y
X = r
```

Thread 3

```
Z = 0
```



- With new Thread 3, printing “foo” **has to be disallowed** (thin-air!)

Thread 1

```

a = X
if (a == 1) {
  Y = a
  print ("foo")
} else {
  b = Z
  Y = b
}

```

Thread 2

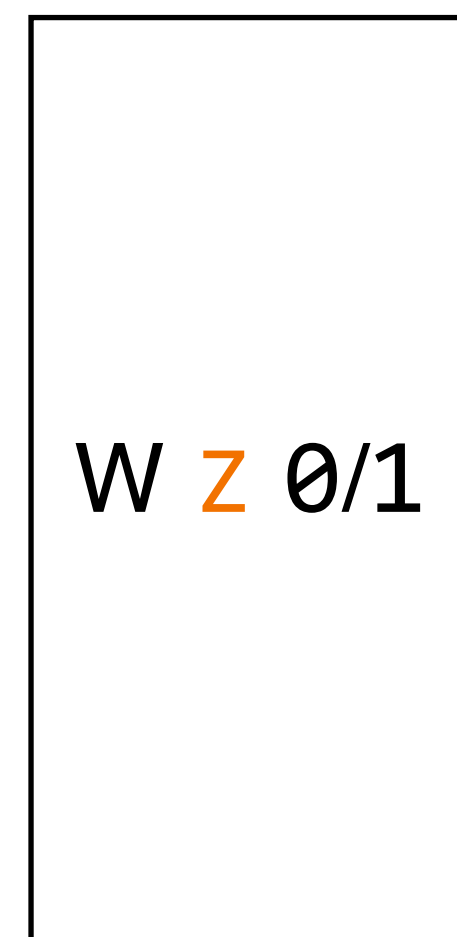
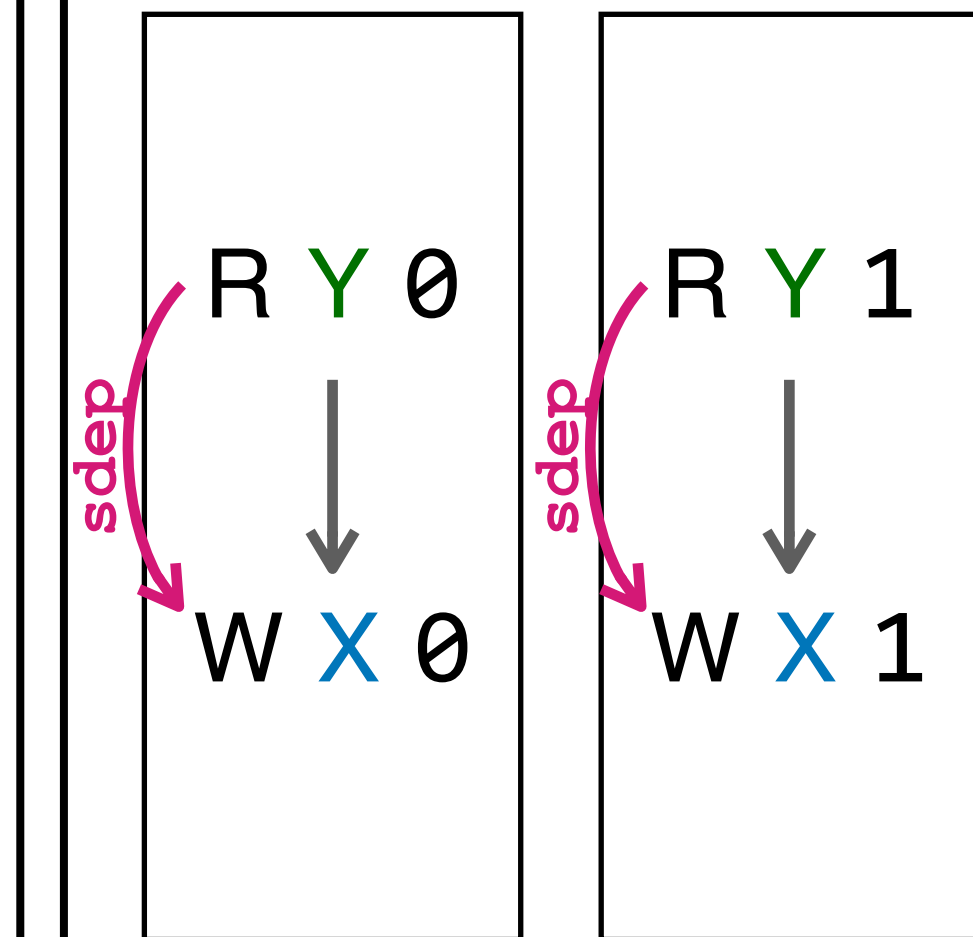
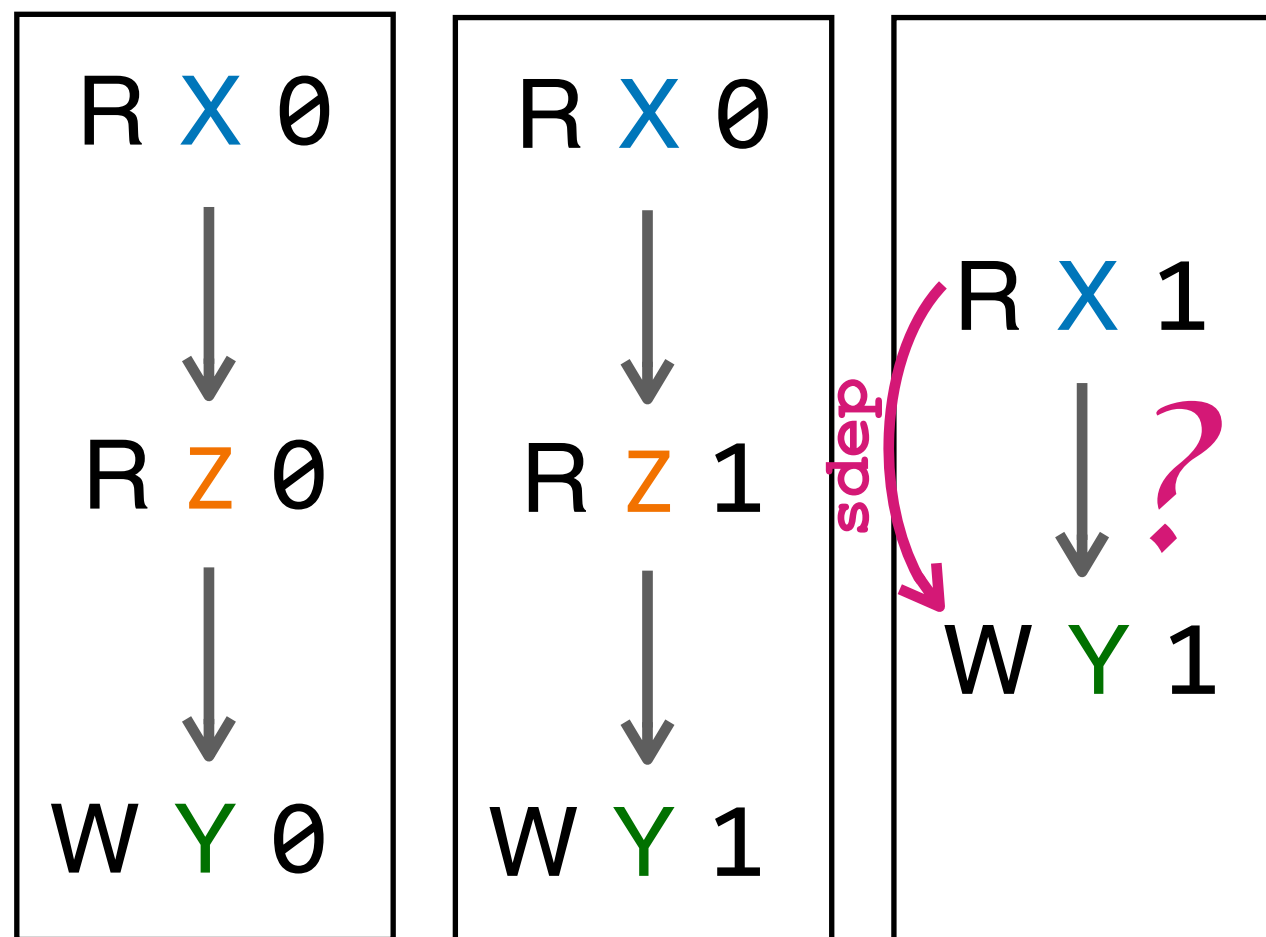
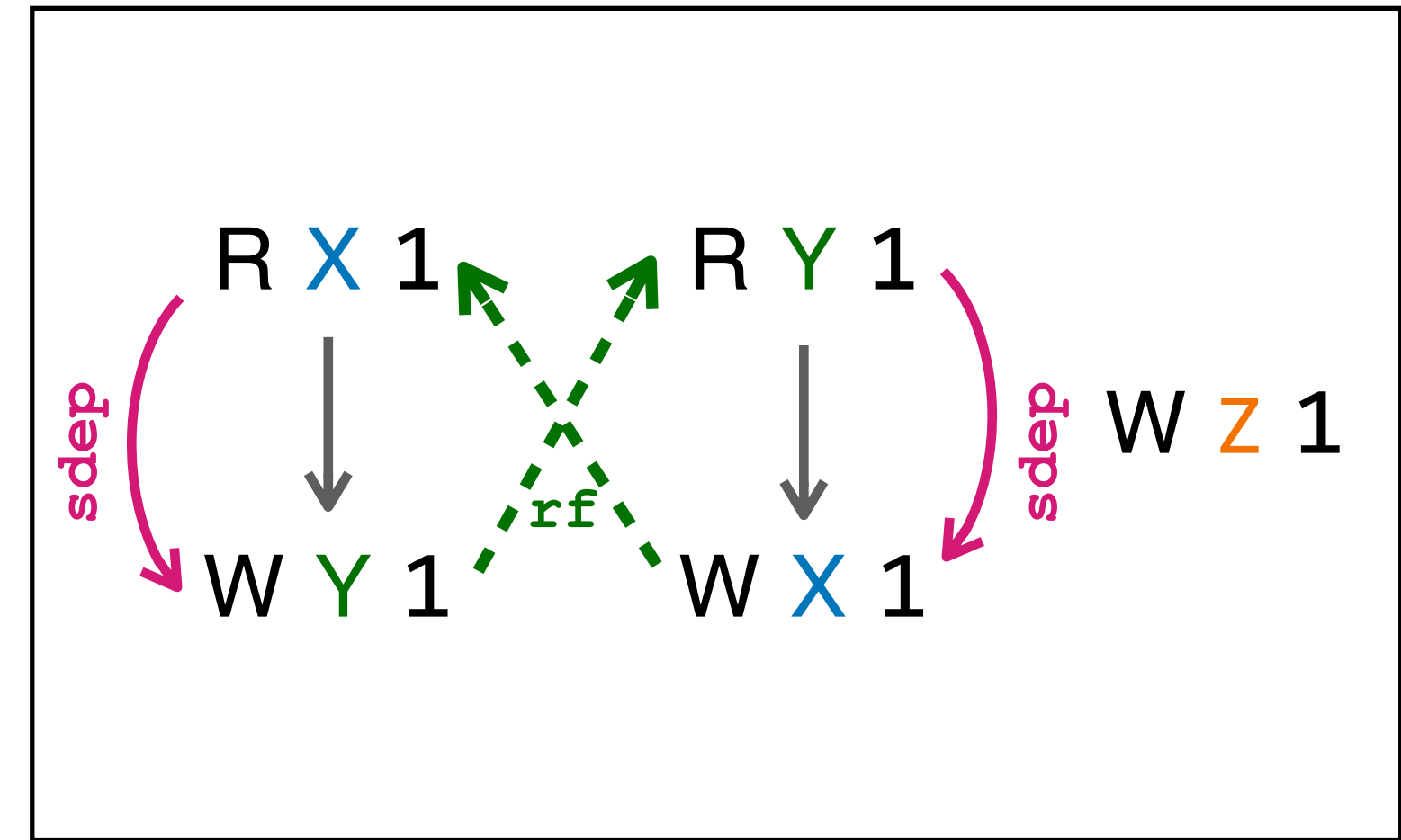
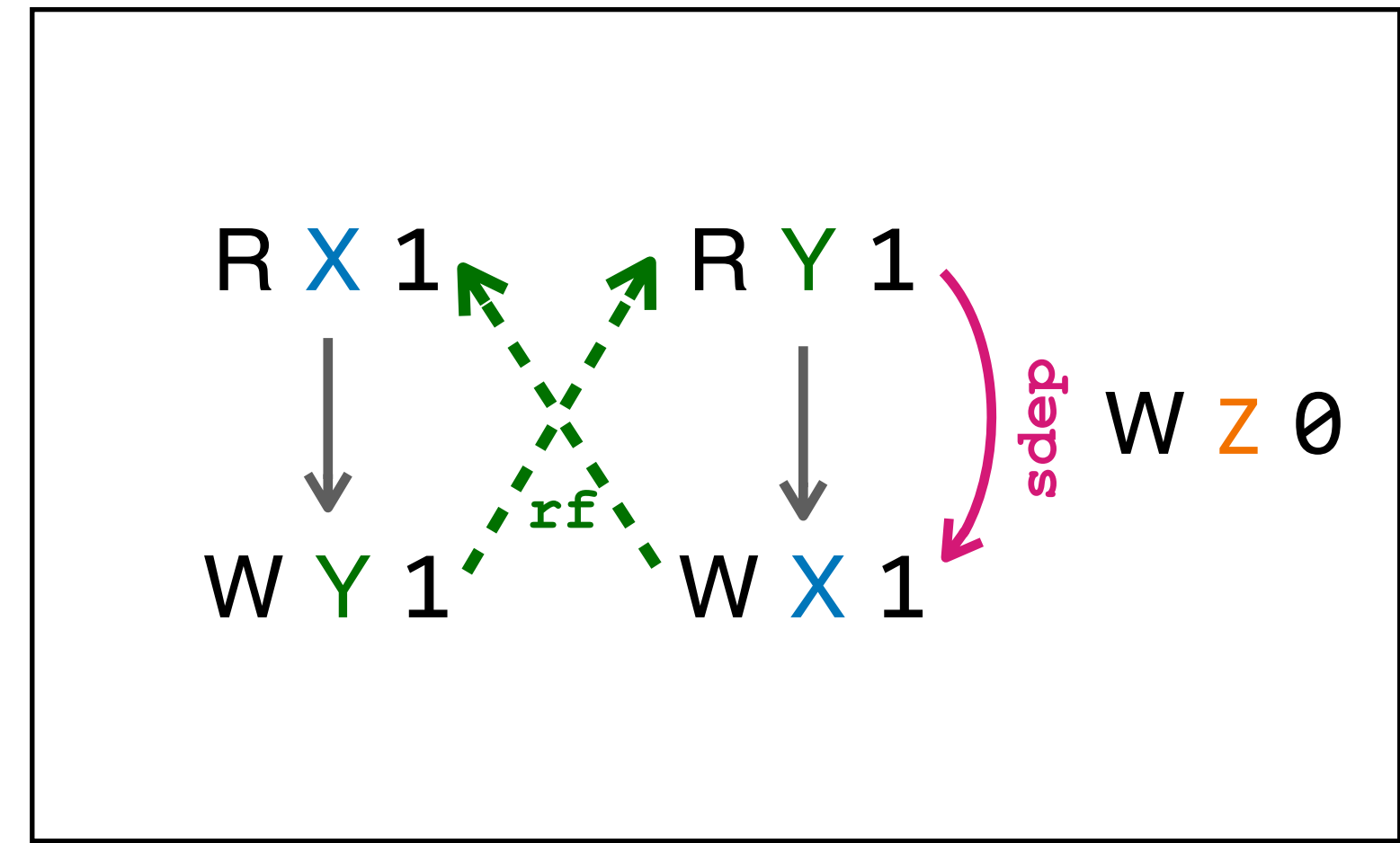
```

r = Y
X = r

```

Thread 3

```
Z = 0/1
```



Step 2 (**sdep** calculation) cannot be thread local!

Thread 1

```
a = X
if (a == 1) {
  Y = a
  print ("foo")
} else {
  b = Z
  Y = b
}
```

Thread 2

```
r = Y
X = r
```

Thread 3

```
Z = your_favorite_litmus_text()
```

```
U = 1
```

```
f = V
```

```
V = f rel
```

```
g = V
```

```
V = 1
```

```
d = FAA(V,1) acq
```

```
e = U
```

```
return(d==1 && e==0 && f==1 && g==2)
```

Step 2 (**sdep** calculation) depends on
Step 3 (the consistency predicate)!

Thread 1

```
a = X
if (a == 1) {
  Y = a
  print ("foo")
} else {
  b = Z
  Y = b
}
```

Thread 2

```
r = Y
X = r
```

Thread 3

```
Z = your_favorite_litmus_test()
```

Substitution of Equivalents — “sanity condition” for weak memory models:

- If $f()$ always returns \emptyset in a memory model M , then $f()$ and \emptyset should be a equivalent in M
(assuming $f()$ uses a disjoint set of locations wrt rest of the program).

Reasoning-aware **sdep**? :(

- **sdep** calculation has to take into account our reasoning principles.
 - No thin air values: $f()$ never returns 1 in some (possibly inconsistent) execution
 \implies **sdep** must exist
 - A new (sound) program logic can prove that $f()$ never returns 1
 \implies **sdep** must exist
- We have a **memory model for reasoning** (weaker than the “real” model).
- For reasoning to be potentially precise, **sdep** needs to take into account the full consistency predicate.

The source of the problem

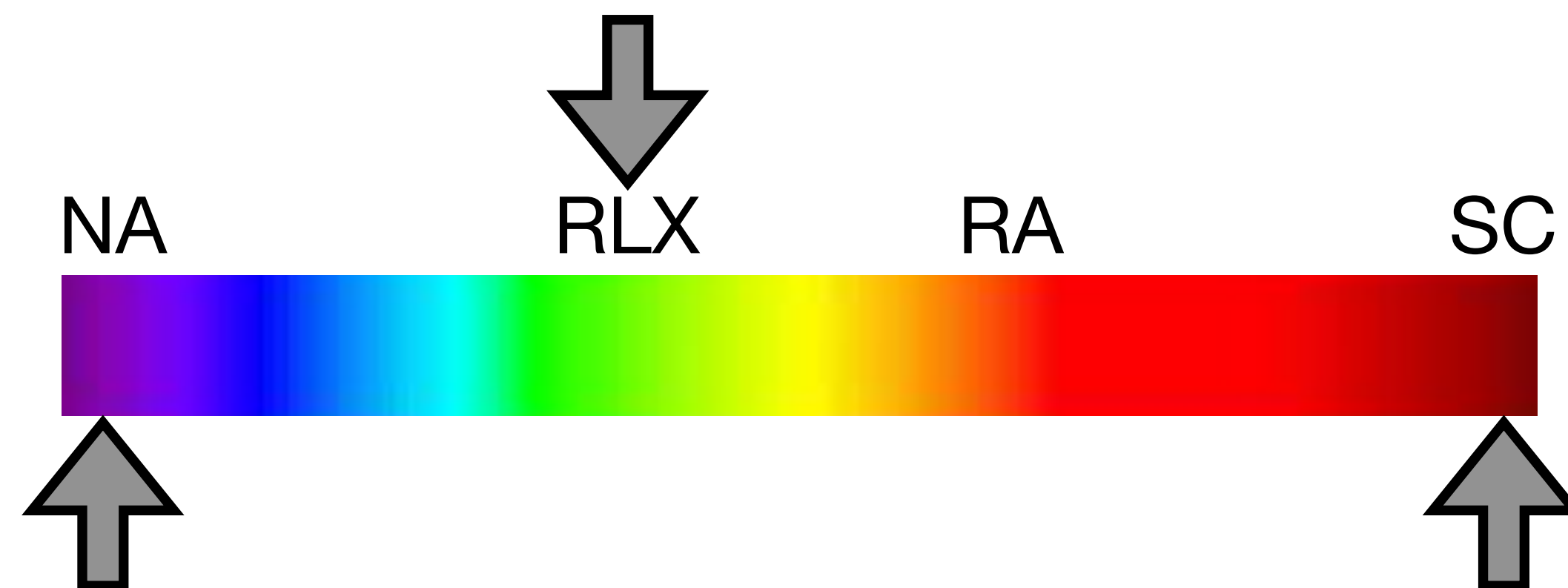
- Semantic dependencies are “**dynamic**” rather than “**static**”:
 - **sdep** \iff the model allows a thread to read some value at a certain program point.
 - Event-structure-based / pomset models / “Promising Semantics” capture such dynamic dependencies.

$$\begin{array}{l} \text{(MACHINE: NORMAL)} \\ \langle \mathcal{T}(\tau), P_G, M \rangle \rightarrow^+ \langle T', P_{G'}, M' \rangle \\ \langle T', P_{G'}, M' \rangle \rightarrow^* \langle _ _ \emptyset, _ _ \rangle \\ \hline \langle \mathcal{T}, P_G, M \rangle \rightarrow \langle \mathcal{T}[\tau \mapsto T'], P_{G'}, M' \rangle \end{array}$$

- The approach we discussed fails to do so.

A fresh look on the out-of-thin-air problem

- The discussion about the OOTA problem in C/C++ revolves around `memory_order_relaxed`
 - *Is it indeed expensive to forbid RW reordering of relaxed accesses?*
 - More provocatively: *do we really need relaxed writes?*



- A (more practical?) challenging problem arises with:
 - **Strong accesses (SC) or mutexes** that allow races
 - **Weak accesses (non-atomic)** that allow optimizations, **including load introduction**