

CS4212: Compiler Design

Week 13:

General Dataflow Analysis Framework

Register Allocation

Wrap-Up

Ilya Sergey

ilya@nus.edu.sg

ilyasergey.net/CS4212/

Previous Lectures: Liveness

- A variable v is *live* at a program point if v is defined before the program point and used after it.
- Liveness is defined in terms of where variables are *defined* and where variables are *used*
- Liveness analysis: Compute the live variables between each statement.
 - May be *conservative* (i.e. it may claim a variable is live when it isn't) so because that's a safe approximation
 - To be useful, it should be more *precise* than simple scoping rules.

Simple Liveness Analysis Algorithm

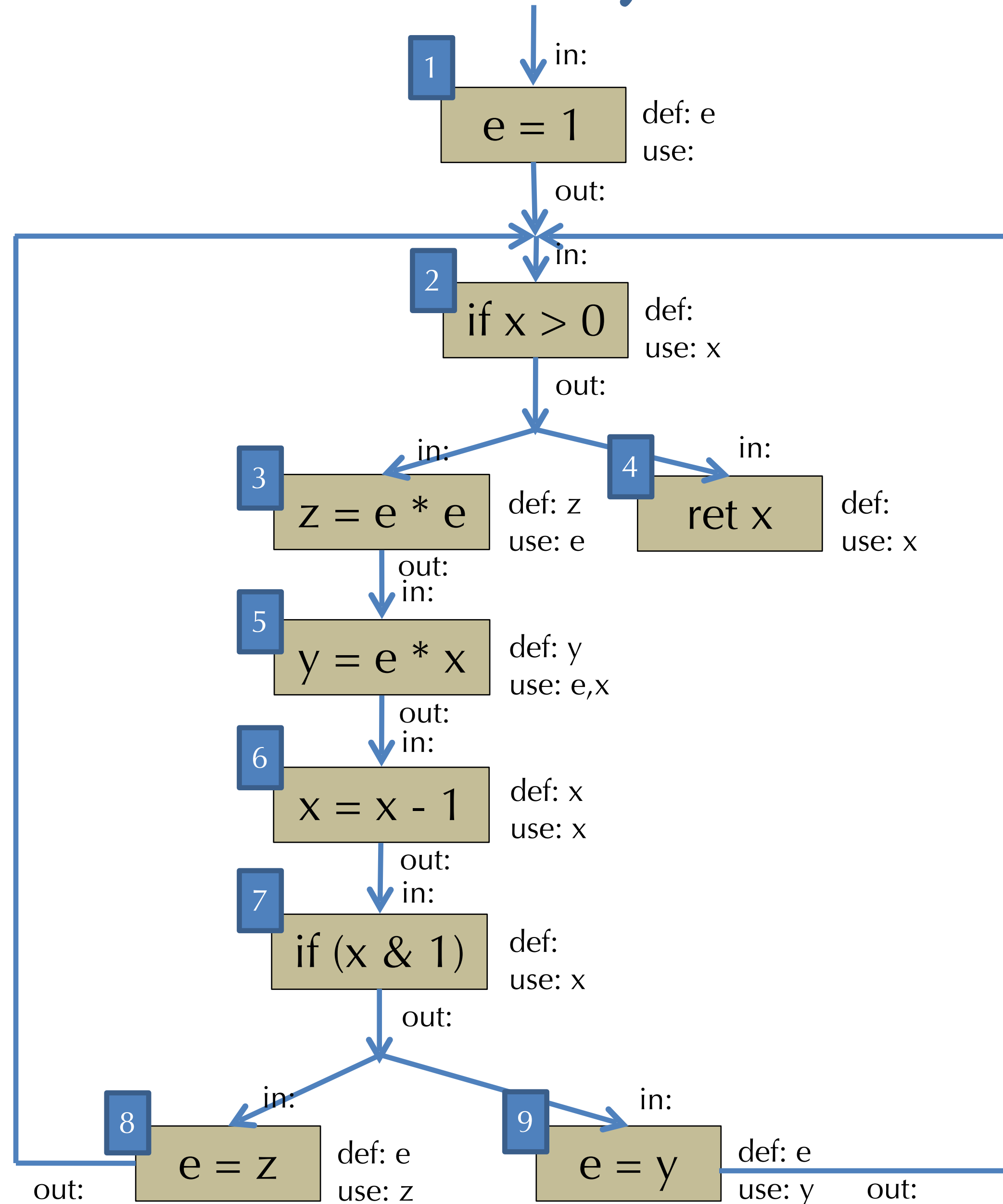
```
for all n, in[n] :=  $\emptyset$ , out[n] :=  $\emptyset$ 
repeat until no change in 'in' and 'out'
  for all n
    out[n] :=  $\bigcup_{n' \in \text{succ}[n]} \text{in}[n']$ 
    in[n] := use[n]  $\cup$  (out[n] - def[n])
  end
end
```

- Finds a *fixpoint* of the **in** and **out** equations.

Example Liveness Analysis

- Example flow graph:

```
e = 1;  
while(x>0) {  
  z = e * e;  
  y = e * x;  
  x = x - 1;  
  if (x & 1) {  
    e = z;  
  } else {  
    e = y;  
  }  
}  
return x;
```



Example Liveness Analysis

Each iteration update:

$$\text{out}[n] := \bigcup_{n' \in \text{succ}[n]} \text{in}[n']$$

$$\text{in}[n] := \text{use}[n] \cup (\text{out}[n] - \text{def}[n])$$

• Iteration 1:

$$\text{in}[2] = x$$

$$\text{in}[3] = e$$

$$\text{in}[4] = x$$

$$\text{in}[5] = e, x$$

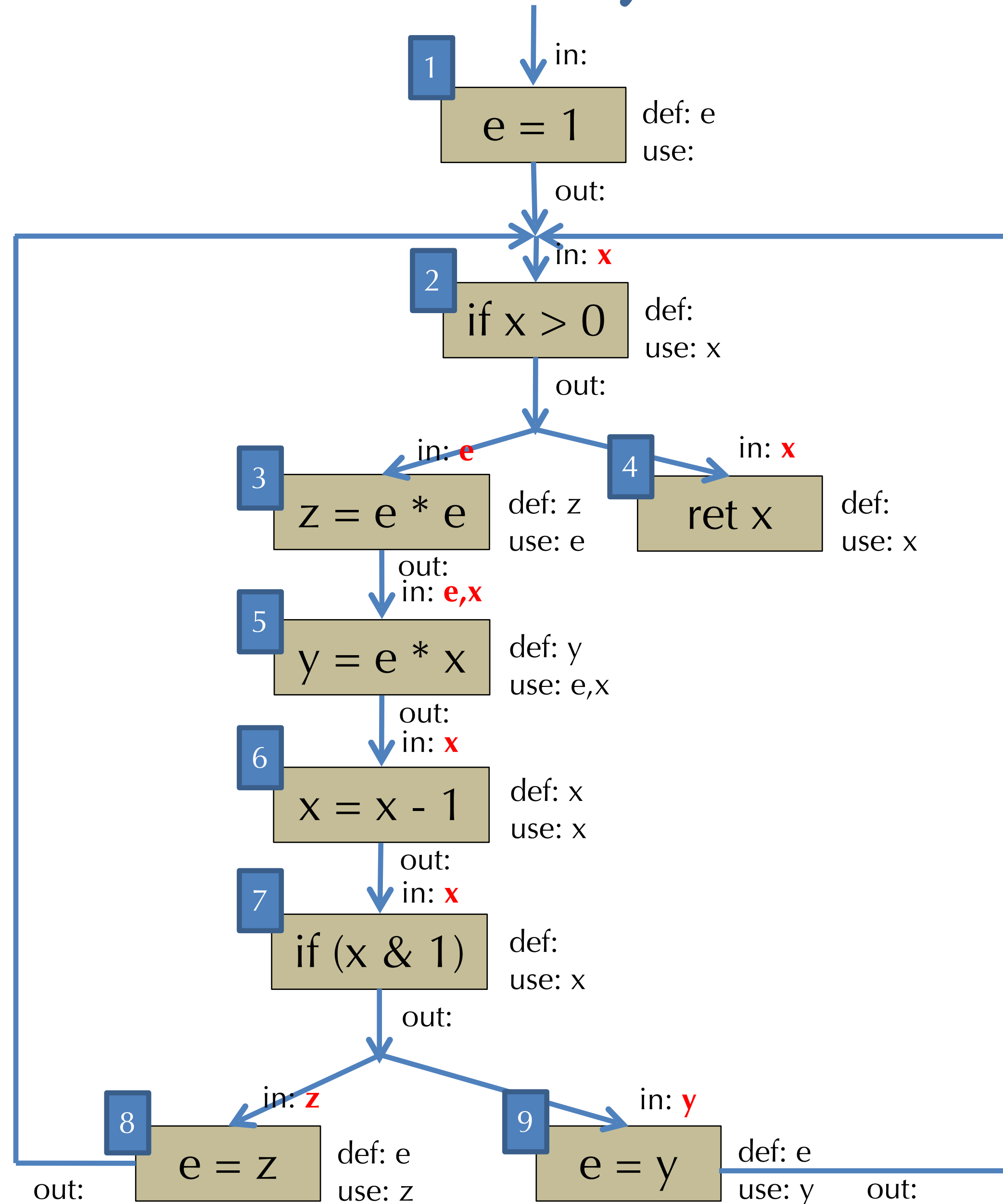
$$\text{in}[6] = x$$

$$\text{in}[7] = x$$

$$\text{in}[8] = z$$

$$\text{in}[9] = y$$

(showing only updates that make a change)



Example Liveness Analysis

Each iteration update:

$$\text{out}[n] := \bigcup_{n' \in \text{succ}[n]} \text{in}[n']$$

$$\text{in}[n] := \text{use}[n] \cup (\text{out}[n] - \text{def}[n])$$

- Iteration 2:

out[1] = x

in[1] = x

out[2] = e, x

in[2] = e, x

out[3] = e, x

in[3] = e, x

out[5] = x

out[6] = x

out[7] = z, y

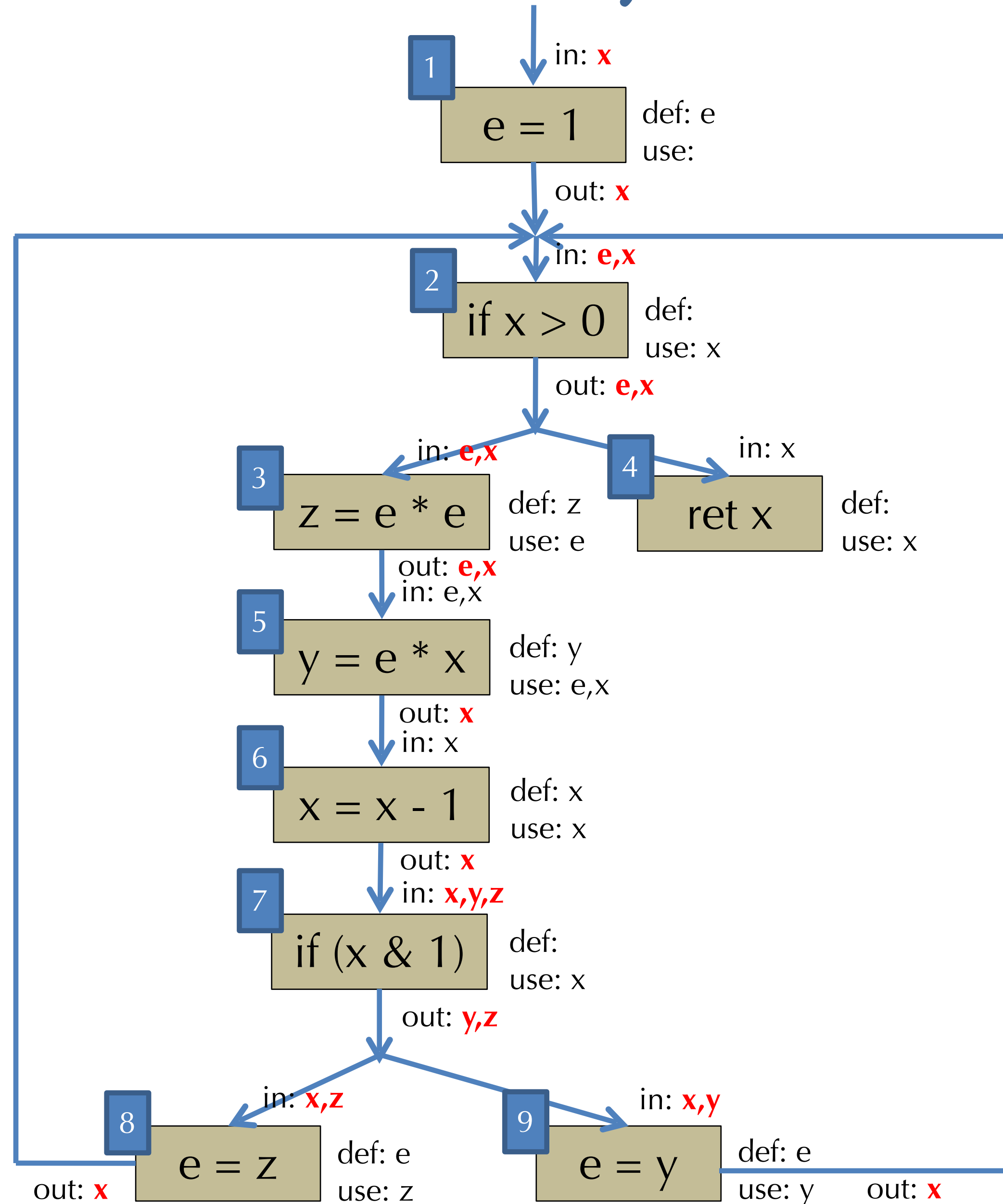
in[7] = x, z, y

out[8] = x

in[8] = x, z

out[9] = x

in[9] = x, y



Example Liveness Analysis

Each iteration update:

$$\text{out}[n] := \bigcup_{n' \in \text{succ}[n]} \text{in}[n']$$
$$\text{in}[n] := \text{use}[n] \cup (\text{out}[n] - \text{def}[n])$$

- Iteration 3:

$\text{out}[1] = e, x$

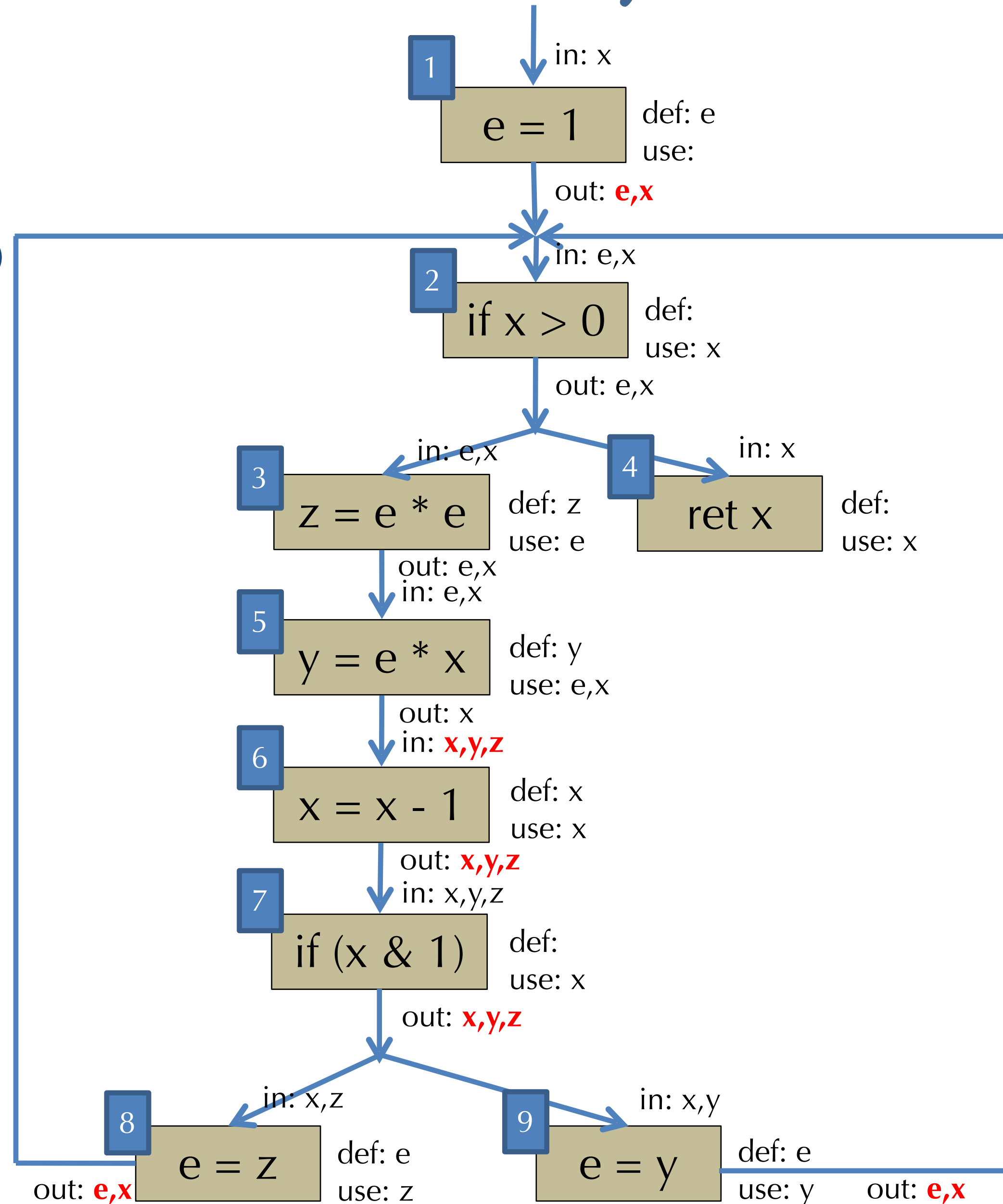
$\text{out}[6] = x, y, z$

$\text{in}[6] = x, y, z$

$\text{out}[7] = x, y, z$

$\text{out}[8] = e, x$

$\text{out}[9] = e, x$



Example Liveness Analysis

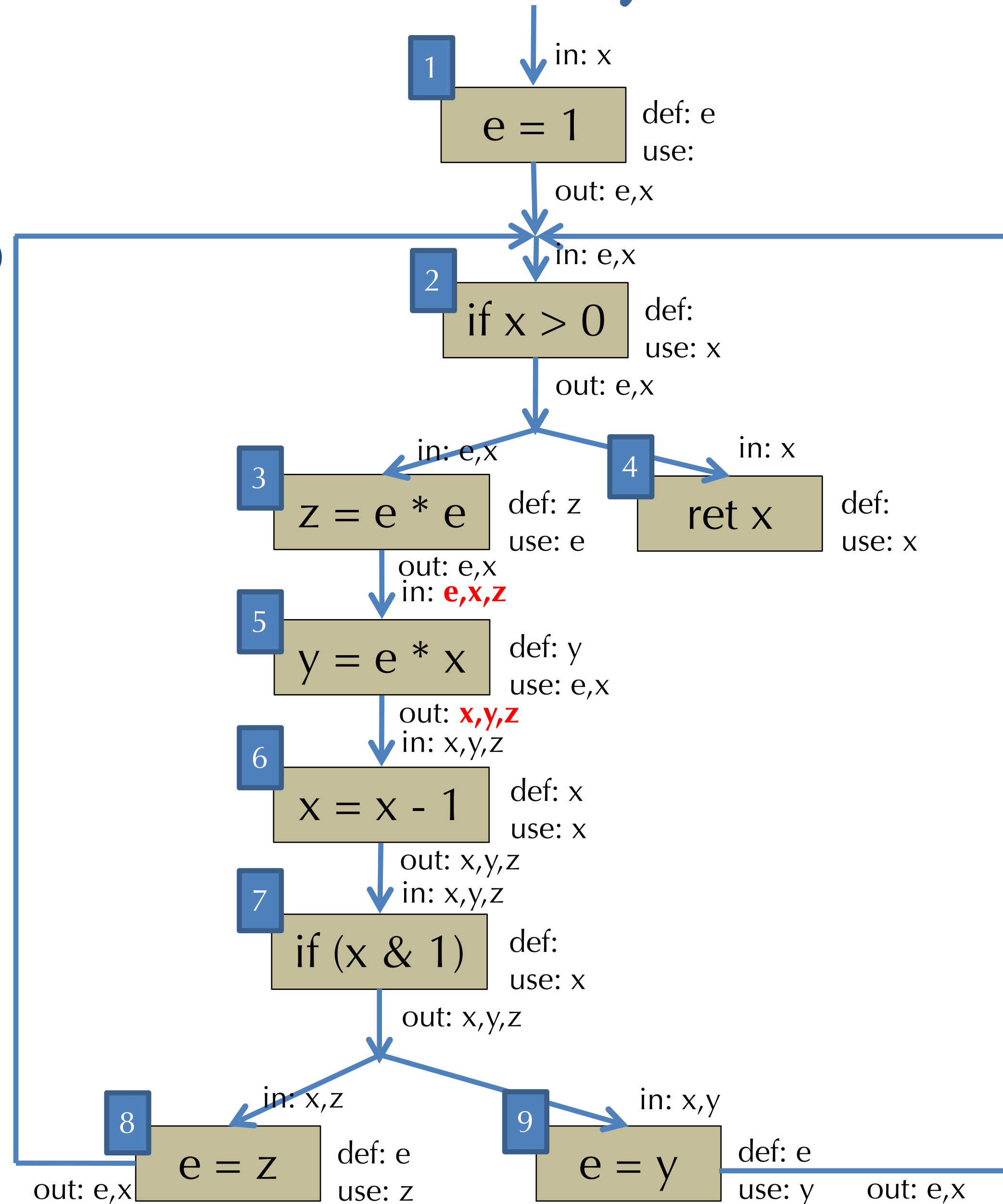
Each iteration update:

$$\text{out}[n] := \bigcup_{n' \in \text{succ}[n]} \text{in}[n']$$
$$\text{in}[n] := \text{use}[n] \cup (\text{out}[n] - \text{def}[n])$$

- Iteration 4:

$\text{out}[5] = x, y, z$

$\text{in}[5] = e, x, z$



Example Liveness Analysis

Each iteration update:

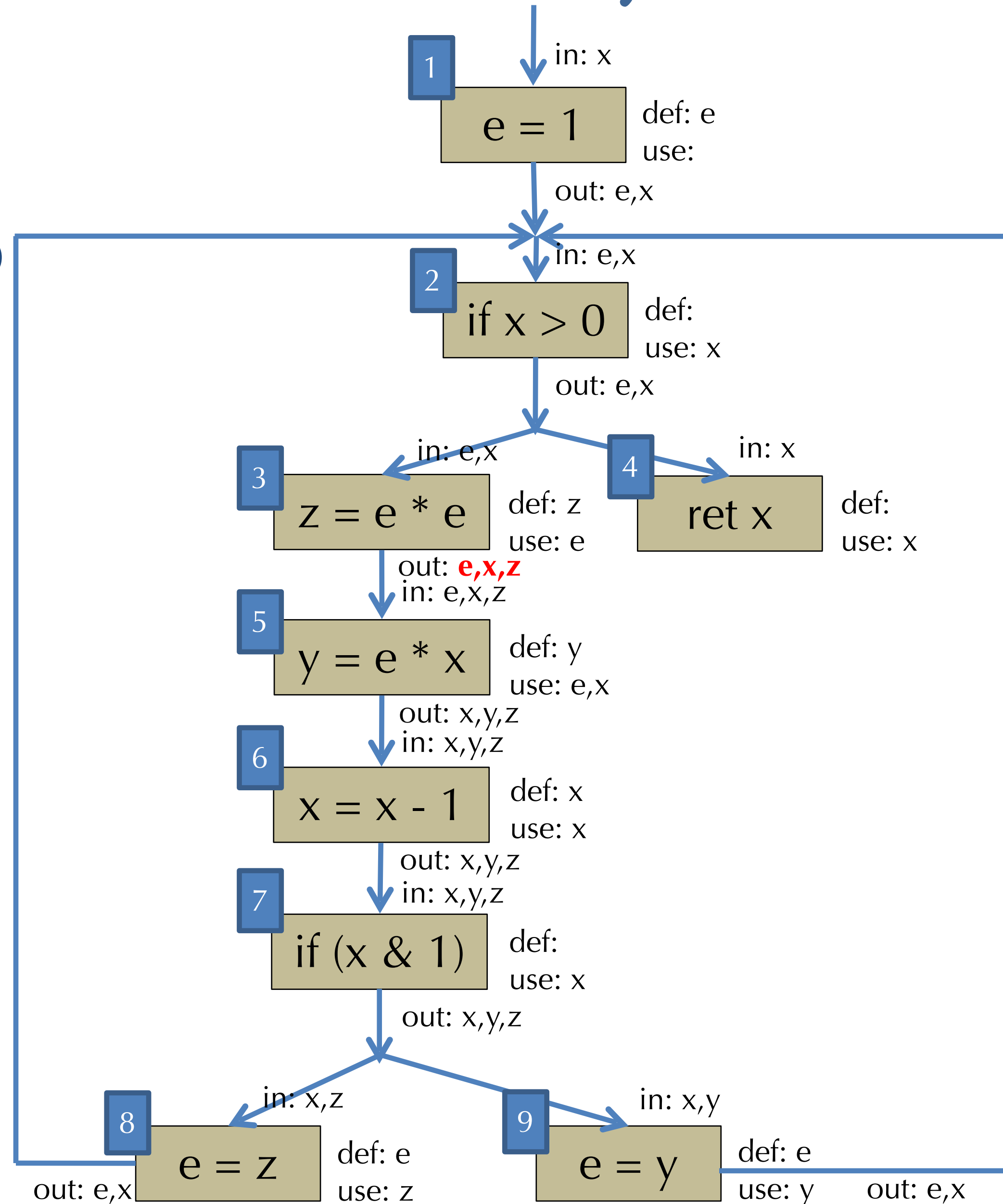
$out[n] := \bigcup_{n' \in succ[n]} in[n']$

$in[n] := use[n] \cup (out[n] - def[n])$

- Iteration 5:

$out[3] = e, x, z$

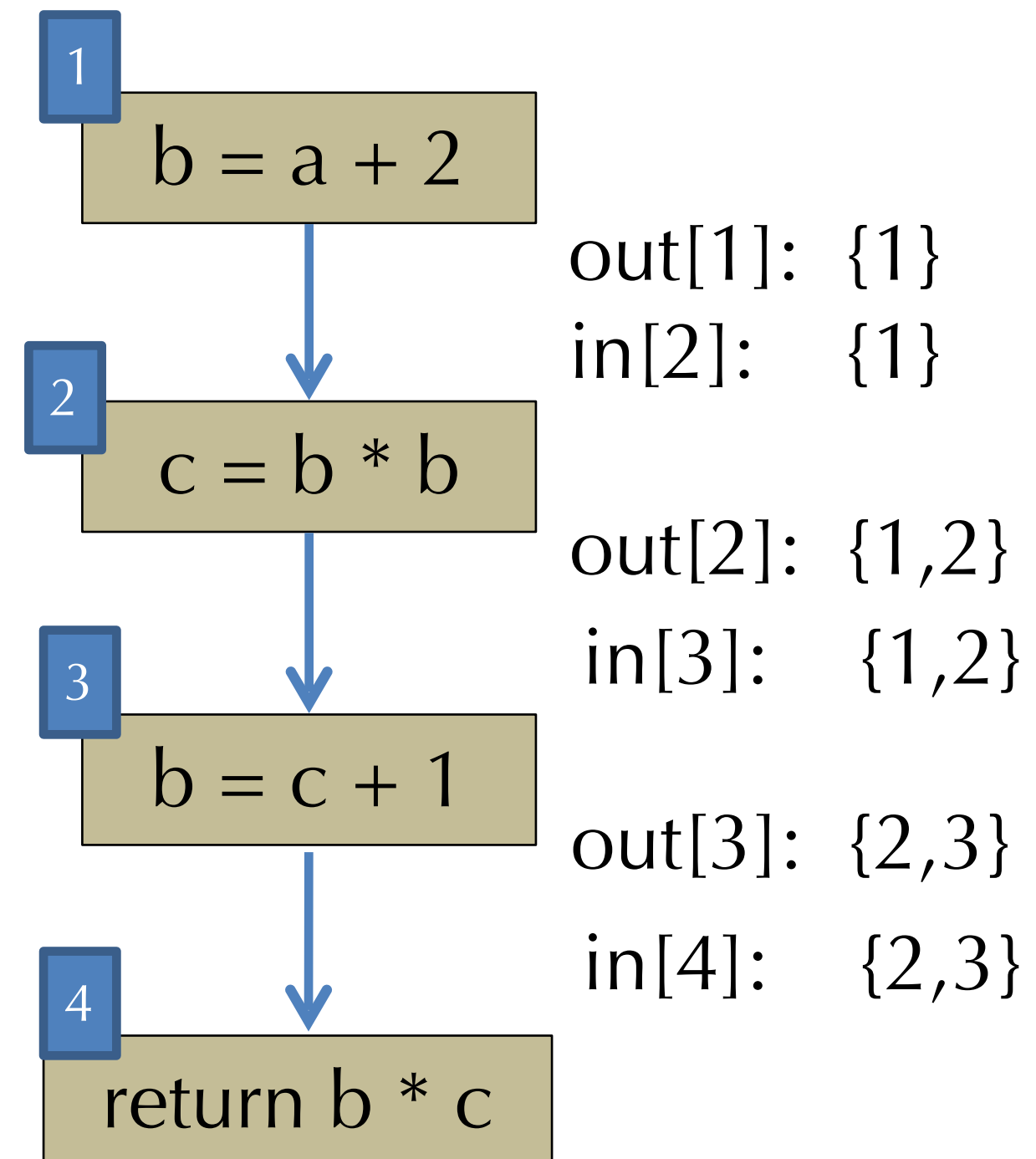
Done!



Reaching Definitions

Example of Reaching Definitions

- Results of computing reaching definitions on this simple CFG:



Reaching Definitions Step 1

- Define the sets of interest for the analysis
 - Let $\text{defs}[a]$ be the set of *nodes* (statements) that define the variable a
- Define $\text{gen}[n]$ and $\text{kill}[n]$ as follows:
 - Quadruple forms n :

	$\text{gen}[n]$	$\text{kill}[n]$
$a = b \text{ op } c$	$\{n\}$	$\text{defs}[a] - \{n\}$
$a = \text{load } b$	$\{n\}$	$\text{defs}[a] - \{n\}$
$\text{store } b, a$	\emptyset	\emptyset
$a = f(b_1, \dots, b_n)$	$\{n\}$	$\text{defs}[a] - \{n\}$
$f(b_1, \dots, b_n)$	\emptyset	\emptyset
$\text{br } L$	\emptyset	\emptyset
$\text{br } a \text{ } L1 \text{ } L2$	\emptyset	\emptyset
$\text{return } a$	\emptyset	\emptyset
- $\text{gen}[n]$ are node's definitions; $\text{kill}[n]$ are the nodes, whose definitions are “shadowed” by n

Reaching Definitions Step 2

- Define the constraints that a reaching definitions solution must satisfy.
- $\text{out}[n] \supseteq \text{gen}[n]$
 - “The definitions that reach the end of a node at least include the definitions generated by the node”
- $\text{in}[n] \supseteq \text{out}[n']$ if n' is in $\text{pred}[n]$
 - “The definitions that reach the beginning of a node include those that reach the exit of its *any* predecessor”
- $\text{out}[n] \cup \text{kill}[n] \supseteq \text{in}[n]$
 - “The definitions that come in to a node either reach the end of the node or are killed by it.”
 - Equivalently: $\text{out}[n] \supseteq \text{in}[n] - \text{kill}[n]$


Reaching Definitions Step 3

- Convert constraints to iterated update equations:
 - $\text{in}[n] := \bigcup_{n' \in \text{pred}[n]} \text{out}[n']$
 - $\text{out}[n] := \text{gen}[n] \cup (\text{in}[n] - \text{kill}[n])$
- Algorithm: initialise $\text{in}[n]$ and $\text{out}[n]$ to \emptyset
 - Iterate the update equations until a fixed point is reached
 - Why does it terminate?
- The algorithm terminates because $\text{in}[n]$ and $\text{out}[n]$ increase only *monotonically*
 - At most to a maximum set that includes all variable definitions in the program
- The algorithm is *precise* because it finds the *smallest* sets that satisfy the constraints.

Available Expressions

Available Expressions

- Idea: want to perform common subexpression elimination:

– $a = x + 1$
...
 $b = x + 1$  $a = x + 1$
...
 $b = a$

- When is it safe?

- This transformation is safe if $x+1$ means computes the same value at both places (i.e., x hasn't been assigned).

- “ $x+1$ ” is an *available expression*

- Dataflow values:

- $\text{in}[n]$ = set of *nodes* whose values are available on entry to n
- $\text{out}[n]$ = set of *nodes* whose values are available on exit of n

Available Expressions Step 1

- Define the sets of values
 - Let $uses[a]$ be the set of *nodes* that use the variable a in their expressions
- Define $gen[n]$ and $kill[n]$ as follows:

Quadruple forms n :	$gen[n]$	$kill[n]$
$a = b \text{ op } c$	$\{n\} - kill[n]$	$uses[a]$
$a = \text{load } b$	$\{n\} - kill[n]$	$uses[a]$
$\text{store } b, a$	\emptyset	$uses[[x]]$ (for all x that may equal a)
$\text{br } L$	\emptyset	\emptyset
$\text{br } a \text{ } L1 \text{ } L2$	\emptyset	\emptyset
$a = f(b_1, \dots, b_n)$	\emptyset	$uses[a] \cup uses[[x]]$ (for all x)
$f(b_1, \dots, b_n)$	\emptyset	$uses[[x]]$ (for all x)
$\text{return } a$	\emptyset	\emptyset

Note the need for
“may alias” information...

Note that functions are
assumed to be impure...

- $gen[n]$ — node itself represents new available expression
- $kill[n]$ — nodes whose expressions no longer available after n

Available Expressions Step 2

- Define the constraints that an available expressions solution must satisfy.
- $out[n] \supseteq gen[n]$
 - “The expressions made available by n that reach the end of the node”
- $in[n] \subseteq out[n']$ if n' is in $pred[n]$
 - “The expressions available at the beginning of a node include those that reach the exit of every predecessor”
- $out[n] \cup kill[n] \supseteq in[n]$
 - “The expressions available on entry either reach the end of the node or are killed by it.”
 - Equivalently: $out[n] \supseteq in[n] - kill[n]$

Note similarities and differences with constraints for “reaching definitions”.

Available Expressions Step 3

- Convert constraints to iterated update equations:
 - $\text{in}[n] := \bigcap_{n' \in \text{pred}[n]} \text{out}[n']$
 - $\text{out}[n] := \text{gen}[n] \cup (\text{in}[n] - \text{kill}[n])$
- Unlike previous algorithms, this one is “shrinking” the set of desired facts
- Algorithm: initialise $\text{in}[n]$ and $\text{out}[n]$ to {set of all nodes}
 - Iterate the update equations until a fixed point is reached
 - Why does the algorithm terminate?
- The algorithm terminates because $\text{in}[n]$ and $\text{out}[n]$ decrease only *monotonically*
 - At most to a minimum of the empty set
- The algorithm is precise because it finds the *largest* sets that satisfy the constraints.

General Dataflow Analysis Framework

Comparing Dataflow Analyses

- Look at the update equations in the inner loop of the analyses
- Liveness:
 - Let $\text{gen}[n] = \text{use}[n]$ and $\text{kill}[n] = \text{def}[n]$
 - $\text{out}[n] := \bigcup_{n' \in \text{succ}[n]} \text{in}[n']$ (backward)
 - $\text{in}[n] := \text{gen}[n] \cup (\text{out}[n] - \text{kill}[n])$
- Reaching Definitions:
 - $\text{in}[n] := \bigcup_{n' \in \text{pred}[n]} \text{out}[n']$ (forward)
 - $\text{out}[n] := \text{gen}[n] \cup (\text{in}[n] - \text{kill}[n])$
- Available Expressions:
 - $\text{in}[n] := \bigcap_{n' \in \text{pred}[n]} \text{out}[n']$ (forward)
 - $\text{out}[n] := \text{gen}[n] \cup (\text{in}[n] - \text{kill}[n])$

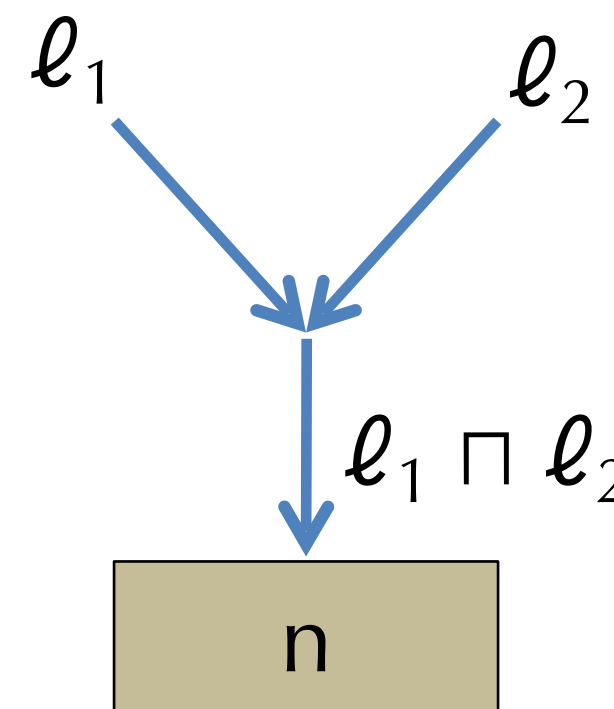
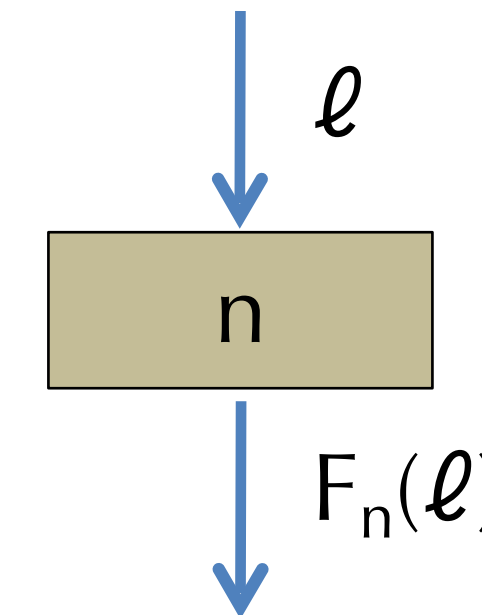
Common Features

- All of these analyses have a *domain* over which they solve constraints.
 - Liveness, the domain is *sets of variables*
 - Reaching defs., Available exprs. the domain is *sets of nodes*
- Each analysis has a notion of *gen[n]* and *kill[n]*
 - Used to explain how information *propagates* across a node: what is added, what is removed.
- Each analysis is propagates information either *forward* or *backward*
 - Forward: *in[n]* defined in terms of predecessor nodes' *out[]*
 - Backward: *out[n]* defined in terms of successor nodes' *in[]*
- Each analysis has a way of aggregating (combining) information from in/out flow
 - Liveness & reaching definitions take union (\cup)
 - Available expressions uses intersection (\cap)
 - Union expresses a property that holds for *some* path (existential)
 - Intersection expresses a property that holds for *all* paths (universal)

(Forward) Dataflow Analysis Framework

A forward dataflow analysis can be characterized by:

1. A domain of dataflow values \mathcal{L}
 - e.g. \mathcal{L} = the powerset of all variables
 - Think of $\ell \in \mathcal{L}$ as a property, then “ $z \in \ell$ ” means “ z has the property”
2. For each node n , a **flow function** $F_n : \mathcal{L} \rightarrow \mathcal{L}$
 - So far we’ve seen $F_n(\ell) = \text{gen}[n] \cup (\ell - \text{kill}[n])$
 - So: $\text{out}[n] = F_n(\text{in}[n])$
 - “If ℓ is a property that holds before the node n , then $F_n(\ell)$ holds after n ”
3. A **combining** operator \sqcap
 - “If we know *either* ℓ_1 *or* ℓ_2 holds on entry to node n , we know at most $\ell_1 \sqcap \ell_2$ ”
 - $\text{in}[n] := \sqcap_{n' \in \text{pred}[n]} \text{out}[n']$



Generic Iterative (Forward) Analysis

```
for all n, in[n] :=  $\top$ , out[n] :=  $\top$ 
repeat until no change
  for all n
    in[n] :=  $\bigcap_{n' \in \text{pred}[n]} \text{out}[n']$ 
    out[n] :=  $F_n(\text{in}[n])$ 
  end
end
```

- Here, $\top \in \mathcal{L}$ (“top”) represents having the “maximum” amount of information.
 - Having “more” information enables more optimizations
 - “Maximum” amount could be inconsistent with the constraints, so we can’t keep it. :-)
 - Iteration refines the answer, eliminating inconsistencies

Structure of \mathcal{L}

- The domain has structure that reflects the “amount” of information for each dataflow value.
- Some dataflow values are more informative than others:
 - Write $\ell_1 \sqsubseteq \ell_2$ whenever ℓ_2 provides at least as much information as ℓ_1 .
 - The dataflow value ℓ_2 is “better” for enabling optimizations.
- Example 1: for available expressions analysis, *larger* sets of nodes are *more informative*.
 - Having a larger set of nodes (equivalently, expressions) available means that there is more opportunity for common subexpression elimination.
 - So: $\ell_1 \sqsubseteq \ell_2$ if and only if $\ell_1 \subseteq \ell_2$
- Example 2: for liveness analysis, *smaller* sets of variables are more informative.
 - Having smaller sets of variables live across an edge means that there are fewer conflicts for register allocation assignments.
 - So: $\ell_1 \sqsubseteq \ell_2$ if and only if $\ell_1 \supseteq \ell_2$

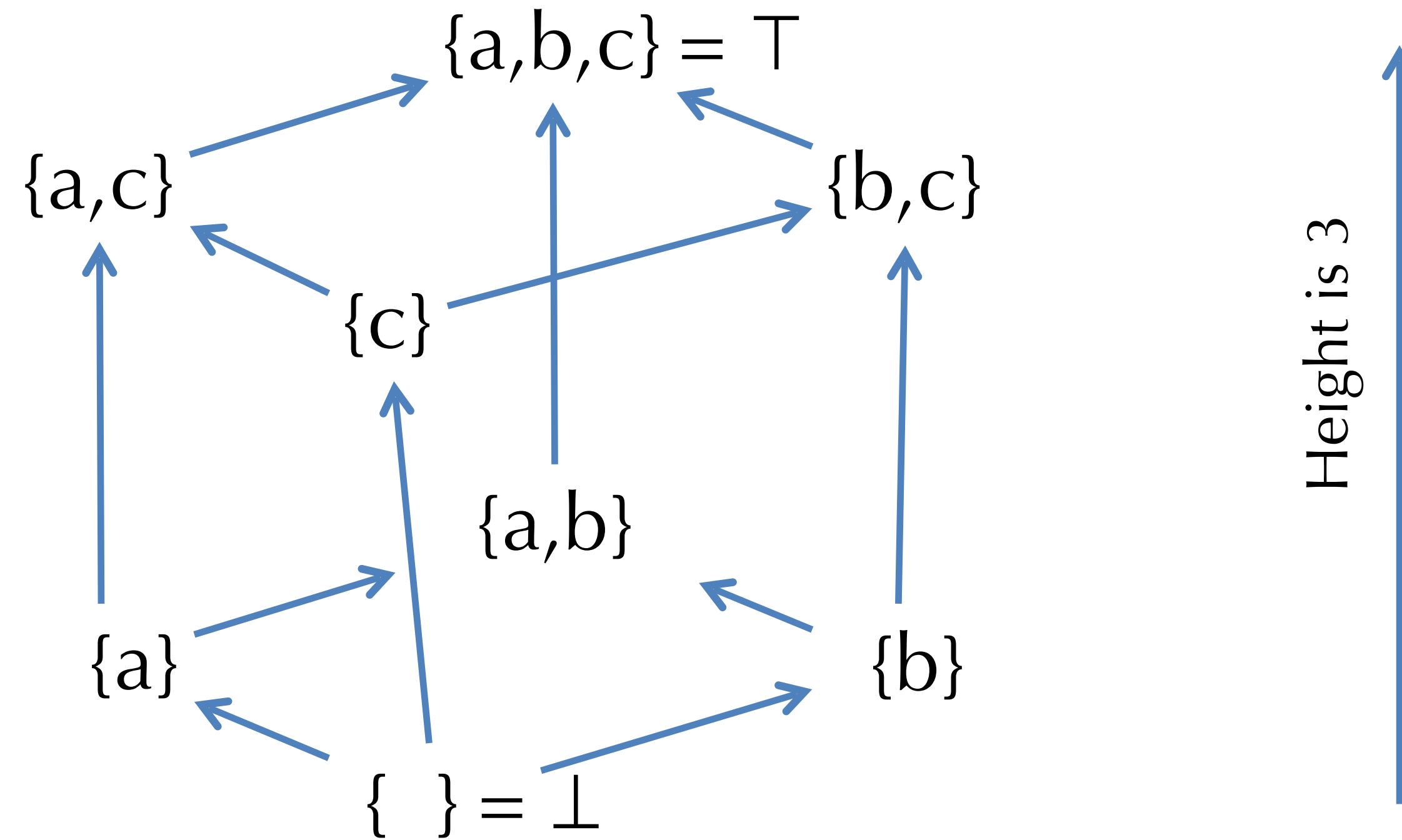
\mathcal{L} as a Partial Order

- \mathcal{L} is a *partial order* defined by the ordering relation \sqsubseteq .
- A partial order is an ordered set.
- Some of the elements might be *incomparable*.
 - That is, there might be $\ell_1, \ell_2 \in \mathcal{L}$ such that neither $\ell_1 \sqsubseteq \ell_2$ nor $\ell_2 \sqsubseteq \ell_1$
- Properties of a partial order:
 - *Reflexivity*: $\ell \sqsubseteq \ell$
 - *Transitivity*: $\ell_1 \sqsubseteq \ell_2$ and $\ell_2 \sqsubseteq \ell_3$ implies $\ell_1 \sqsubseteq \ell_3$
 - *Anti-symmetry*: $\ell_1 \sqsubseteq \ell_2$ and $\ell_2 \sqsubseteq \ell_1$ implies $\ell_1 = \ell_2$
- Examples:
 - Integers ordered by \leq
 - Types ordered by $<$:
 - Sets ordered by \subseteq or \supseteq

Subsets of $\{a,b,c\}$ ordered by \subseteq

$$\begin{array}{c} \ell_2 \\ \uparrow \\ \ell_1 \end{array} \quad \ell_1 \sqsubseteq \ell_2$$

Partial orders are often presented as a Hasse diagram.



order \sqsubseteq is \subseteq

meet \sqcap is \cap

join \sqcup is \cup

Meets and Joins

- The *combining* operator \sqcap is called the “meet” operation.
- It constructs the *greatest lower bound*:
 - $\ell_1 \sqcap \ell_2 \sqsubseteq \ell_1$ and $\ell_1 \sqcap \ell_2 \sqsubseteq \ell_2$
“the meet is a lower bound”
 - If $\ell \sqsubseteq \ell_1$ and $\ell \sqsubseteq \ell_2$ then $\ell \sqsubseteq \ell_1 \sqcap \ell_2$
“there is no greater lower bound”
- Dually, the \sqcup operator is called the “join” operation.
- It constructs the *least upper bound*:
 - $\ell_1 \sqsubseteq \ell_1 \sqcup \ell_2$ and $\ell_2 \sqsubseteq \ell_1 \sqcup \ell_2$
“the join is an upper bound”
 - If $\ell_1 \sqsubseteq \ell$ and $\ell_2 \sqsubseteq \ell$ then $\ell_1 \sqcup \ell_2 \sqsubseteq \ell$
“there is no smaller upper bound”
- A partial order that has all meets and joins is called a *lattice*.
 - If it has just meets, it’s called a *meet semi-lattice*.

Another Way to Describe the (Forward) Algorithm

- Algorithm repeatedly computes (for each node n):
 - $\text{out}[n] := F_n(\text{in}[n])$
- Equivalently: $\text{out}[n] := F_n(\prod_{n' \in \text{pred}[n]} \text{out}[n'])$
 - By definition of $\text{in}[n]$
- We can write this as a simultaneous update of the vector of $\text{out}[n]$ values:
 - Let $x_n = \text{out}[n]$
 - Let $\mathbf{X} = (x_1, x_2, \dots, x_n)$ it's a vector of points in \mathcal{L} corresponding to CFG nodes
 - $\mathbf{F}(\mathbf{X}) = (F_1(\prod_{j \in \text{pred}[1]} \text{out}[j]), F_2(\prod_{j \in \text{pred}[2]} \text{out}[j]), \dots, F_n(\prod_{j \in \text{pred}[n]} \text{out}[j]))$
- Any solution to the constraints is a *fixpoint* \mathbf{X} of \mathbf{F}
 - i.e. $\mathbf{F}(\mathbf{X}) = \mathbf{X}$

Iteration Computes Fixpoints

- Let $\mathbf{X}_0 = (\top, \top, \dots, \top)$
- Each loop through the algorithm apply \mathbf{F} to the old vector:
 $\mathbf{X}_1 = \mathbf{F}(\mathbf{X}_0)$
 $\mathbf{X}_2 = \mathbf{F}(\mathbf{X}_1)$
...
- $\mathbf{F}^{k+1}(\mathbf{X}) = \mathbf{F}(\mathbf{F}^k(\mathbf{X}))$
- A fixpoint is reached when $\mathbf{F}^k(\mathbf{X}) = \mathbf{F}^{k+1}(\mathbf{X})$
 - That's when the algorithm stops.
- Wanted: a *maximal* fixpoint
 - Because that one is more informative/useful for performing optimizations

Monotonicity & Termination

- Each flow function F_n maps lattice elements to lattice elements; to be sensible it should be *monotonic*:
- $F : \mathcal{L} \rightarrow \mathcal{L}$ is *monotonic* iff:
 $\ell_1 \sqsubseteq \ell_2$ implies that $F(\ell_1) \sqsubseteq F(\ell_2)$
 - Intuitively: “If you have more information entering a node, then you have more information leaving the node.”
- Monotonicity lifts point-wise to the function: $\mathbf{F} : \mathcal{L}^n \rightarrow \mathcal{L}^n$
 - vector $(x_1, x_2, \dots, x_n) \sqsubseteq (y_1, y_2, \dots, y_n)$ iff $x_i \sqsubseteq y_i$ for each i
- Note that \mathbf{F} is consistent: $\mathbf{F}(\mathbf{X}_0) \sqsubseteq \mathbf{X}_0$
 - So each iteration moves at least one step down the lattice (for some component of the vector)
 - $\dots \sqsubseteq \mathbf{F}(\mathbf{F}(\mathbf{X}_0)) \sqsubseteq \mathbf{F}(\mathbf{X}_0) \sqsubseteq \mathbf{X}_0$
- Therefore, # steps needed to reach a fixpoint is at most the height H of \mathcal{L} times the number of nodes:
 $O(H_n)$ — height of the lattice

More on Fixpoint Solutions

- Remember constructing LL(1) parse tables

$T \mapsto S\$$
 $S \mapsto ES'$
 $S' \mapsto \varepsilon$
 $S' \mapsto + S$
 $E \mapsto \text{number} \mid (S)$

- $\text{First}(T) = \text{First}(S)$
- $\text{First}(S) = \text{First}(E)$
- $\text{First}(S') = \{ + \}$
- $\text{First}(E) = \{ \text{number}, '(' \}$
- $\text{Follow}(S') = \text{Follow}(S)$
- $\text{Follow}(S) = \{ \$, ')' \} \cup \text{Follow}(S')$

Then: we want the *least* solution to this system of set equations... a *fixpoint* computation. More on these later in the course.

Now: This solution is obtained by starting from taking all First/Follow as \emptyset and then iterating the equations until *fixpoint* is reached.

	number	+	()	\$ (EOF)
T	$\mapsto S\$$		$\mapsto S\$$		
S	$\mapsto E S'$		$\mapsto E S'$		
S'		$\mapsto + S$		$\mapsto \varepsilon$	$\mapsto \varepsilon$
E	$\mapsto \text{num.}$		$\mapsto (S)$		

Dataflow Analysis: Summary

- Many dataflow analyses fit into a common framework.
- Key idea: *iterative solution* of a system of equations (constraints) over a *lattice* of *facts*.
 - Iteration terminates if flow functions are *monotonic*.
 - Solution is obtained as the *greatest* fixpoint is reached via the *meet* operation (\sqcap).
- In the literature, sometimes the definition of the analysis lattice is *reversed*:
 - The most useful/precise information is represented by the bottom element (\perp)
 - Solution is obtained as the *least* fixpoint via iterative application of *join* operator (\sqcup)
 - The two definitions are equivalent modulo the (semi-)lattice *direction*.

Implementation

- See HW6
 - Generic analysis is to be defined in **solver.ml**
 - Control-Flow Graphs are defined in **cfg.ml**
- Analysis example: **liveness.ml**
- Printing analysis results, e.g., liveness:

`./printanalysis.native -live llprograms/analysis2.ll`

Let's have a short break

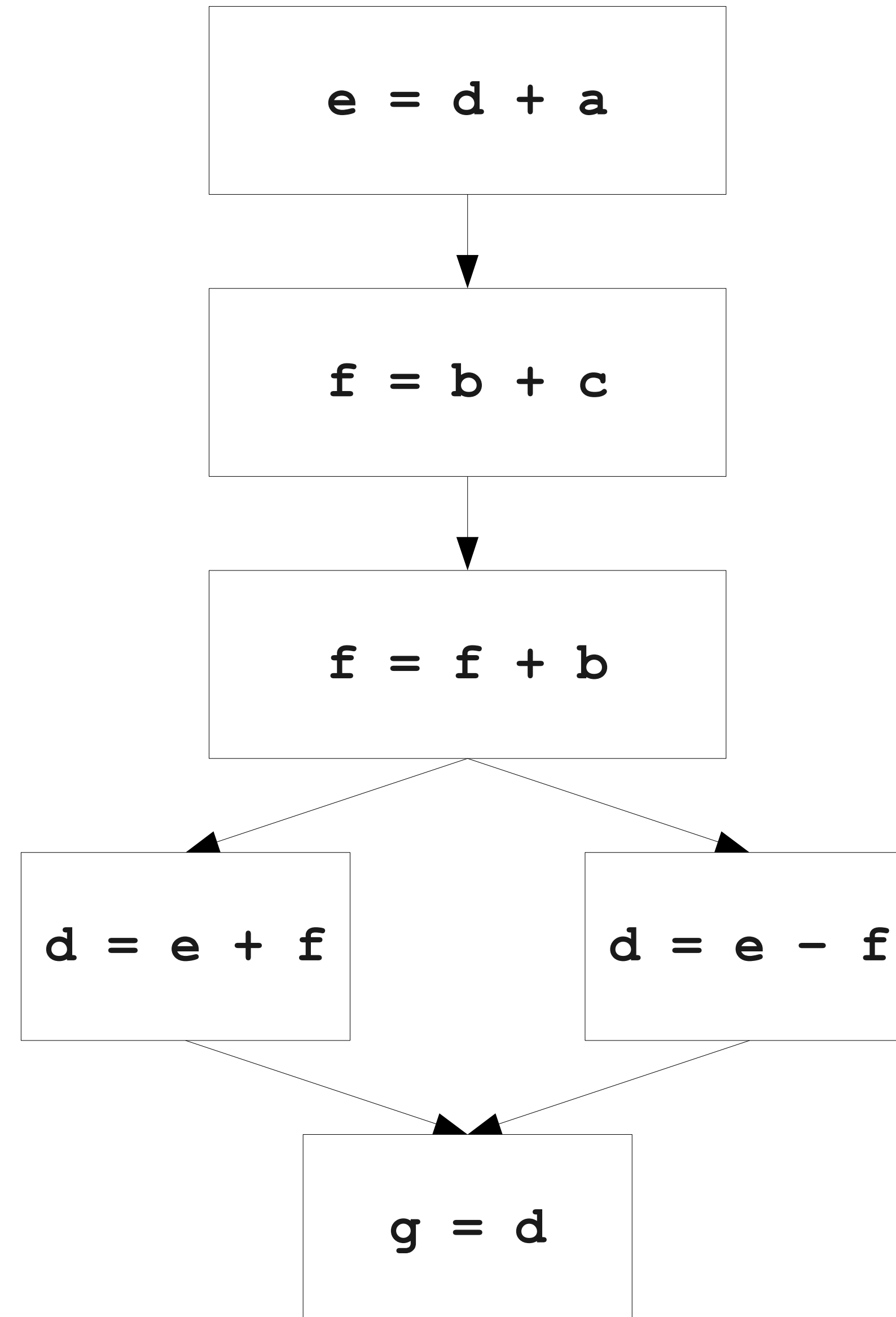
Register Allocation

Register Allocation Problem

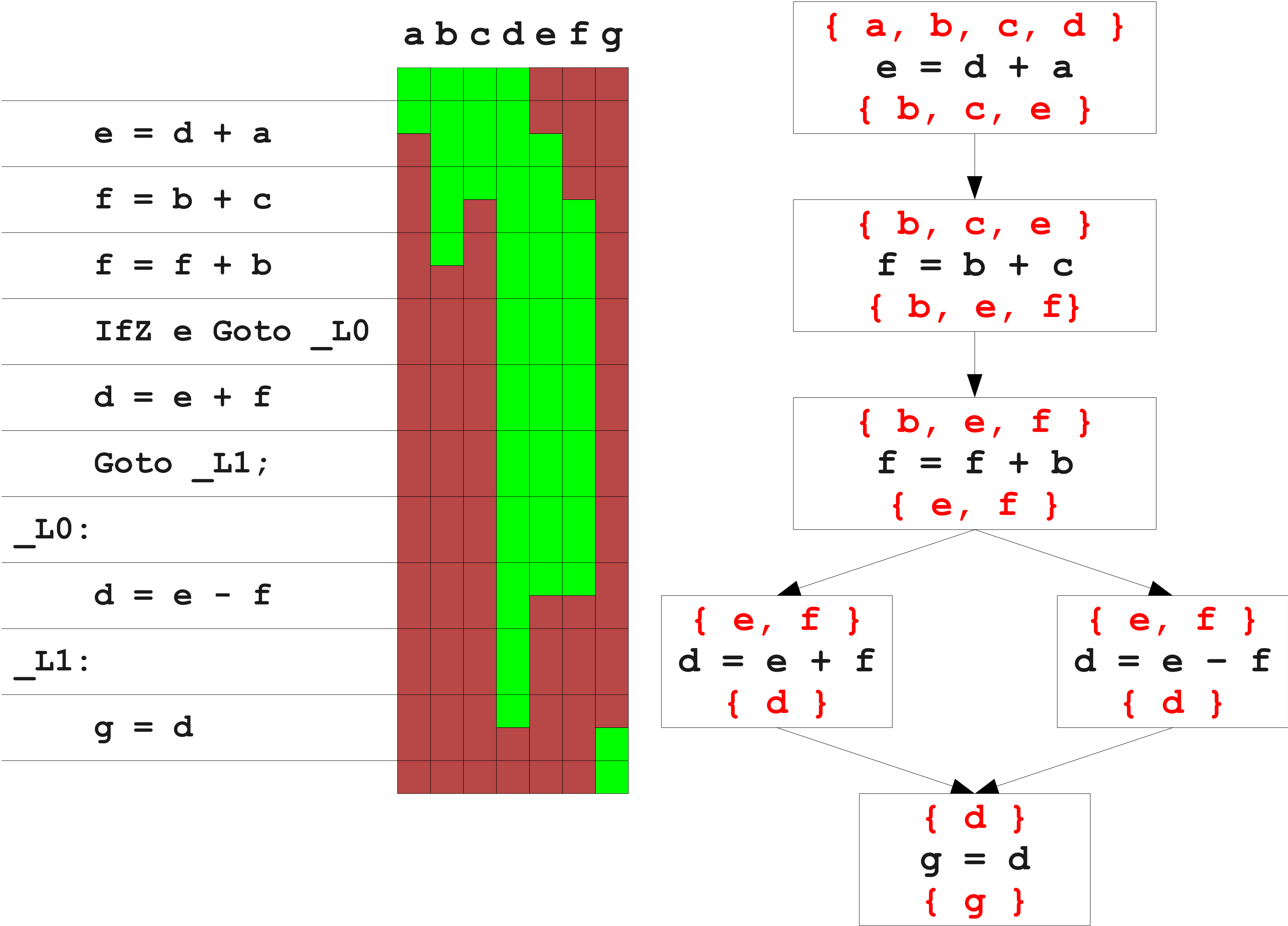
- Given: an IR program that uses an unbounded number of temporaries
 - e.g. the uids of our LLVM programs
- Find: a mapping from temporaries to machine registers such that
 - program semantics is preserved (i.e. the behaviour is the same)
 - register usage is maximised
 - moves between registers are minimised
 - calling conventions / architecture requirements are obeyed
- *Stack Spilling*
 - If there are k registers available and $m > k$ temporaries are *live* at the same time, then not all of them will fit into registers.
 - So: “spill” the excess temporaries to the stack.

Linear-Scan Register Allocation

```
e = d + a
f = b + c
f = f + b
IfZ e Goto _L0
d = e + f
Goto _L1;
_L0:
  d = e - f
_L1:
  g = d
```



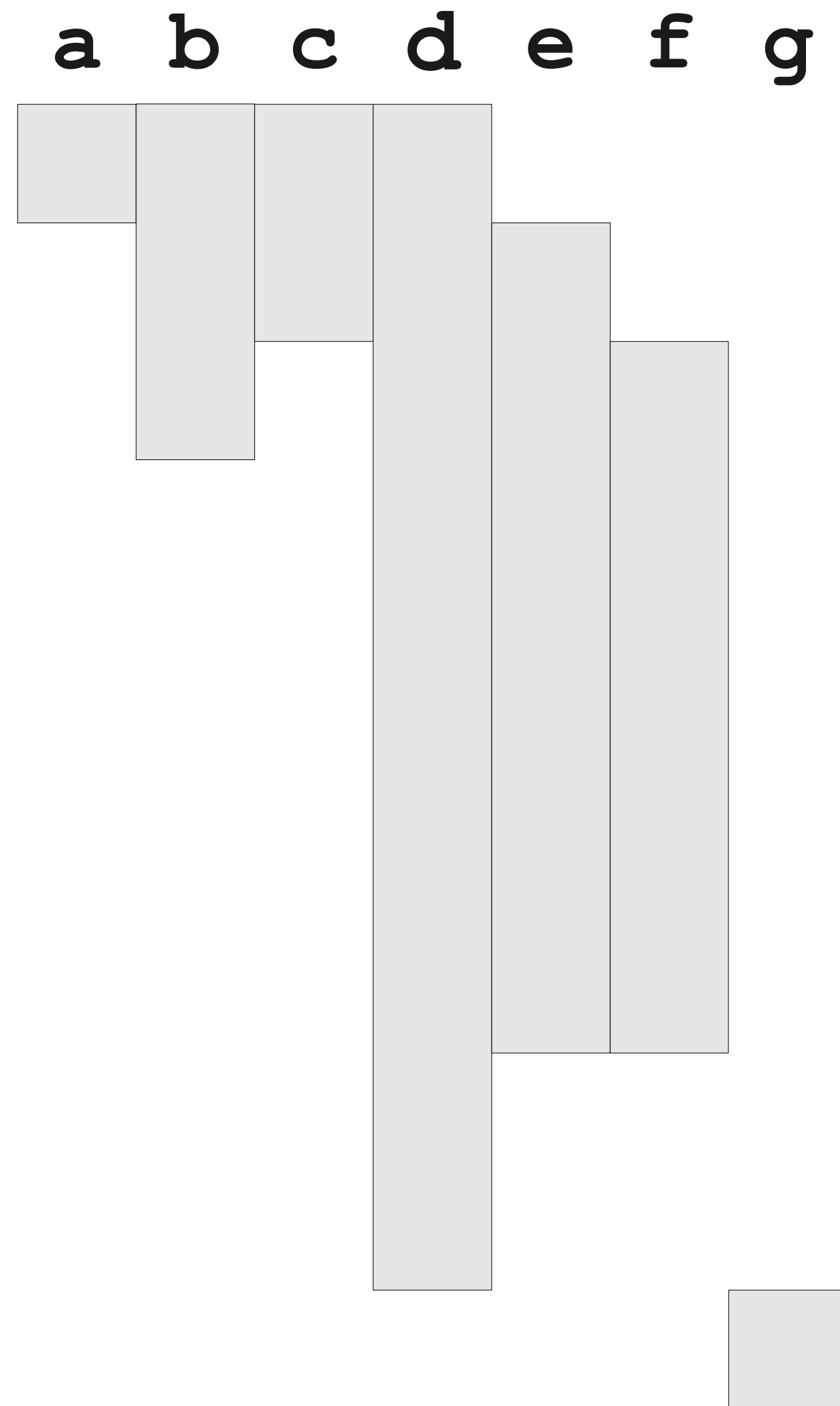
Linear-Scan Register Allocation



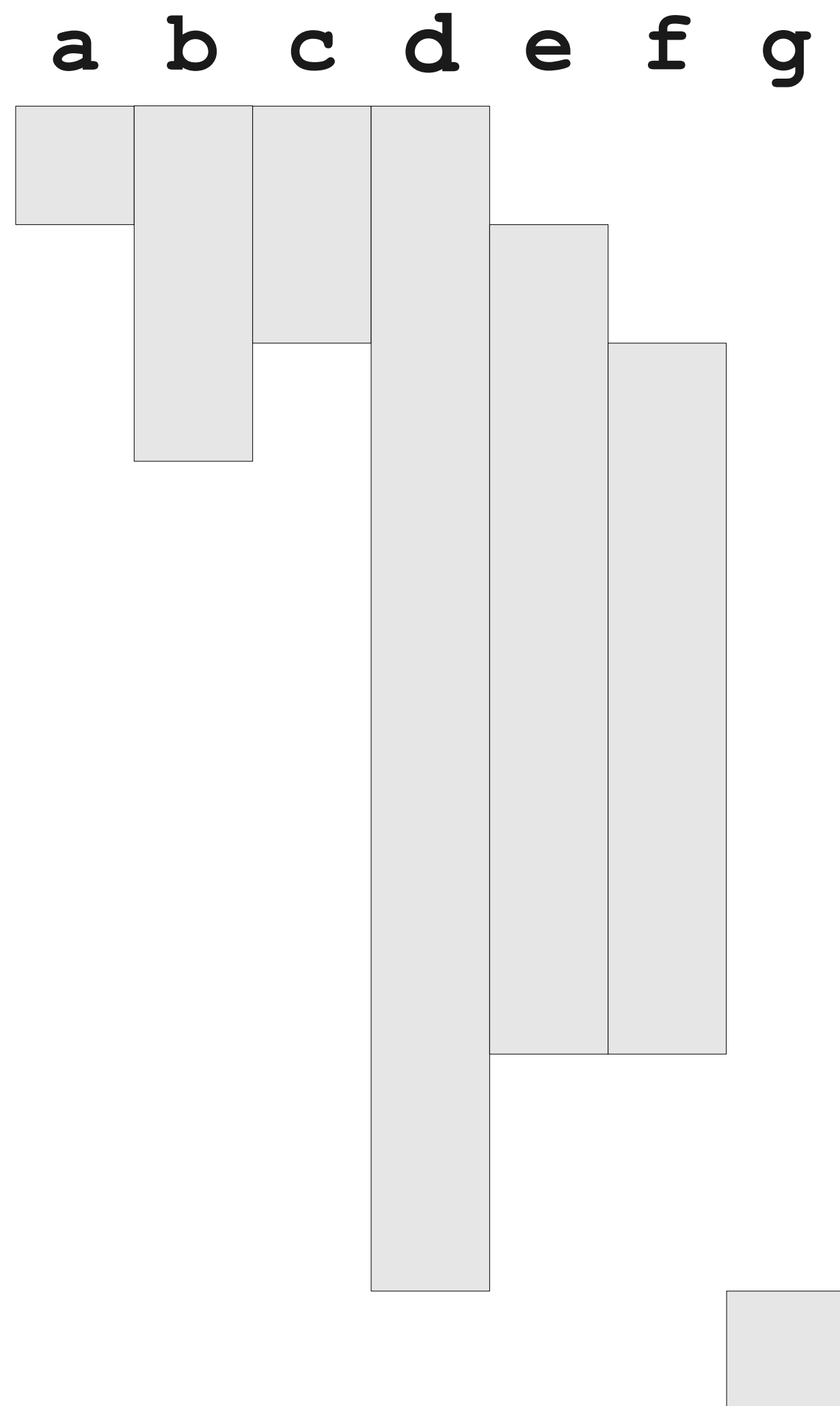
Linear-Scan Register Allocation

Idea: sweep the program *top-down*,
allocating registers for *live* variables and *evicting* non-live ones.

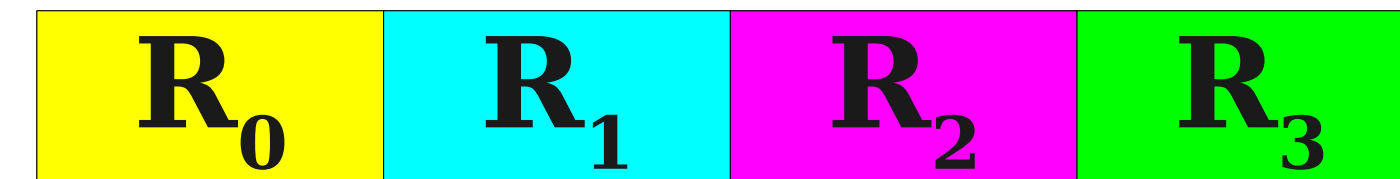
Linear-Scan Register Allocation



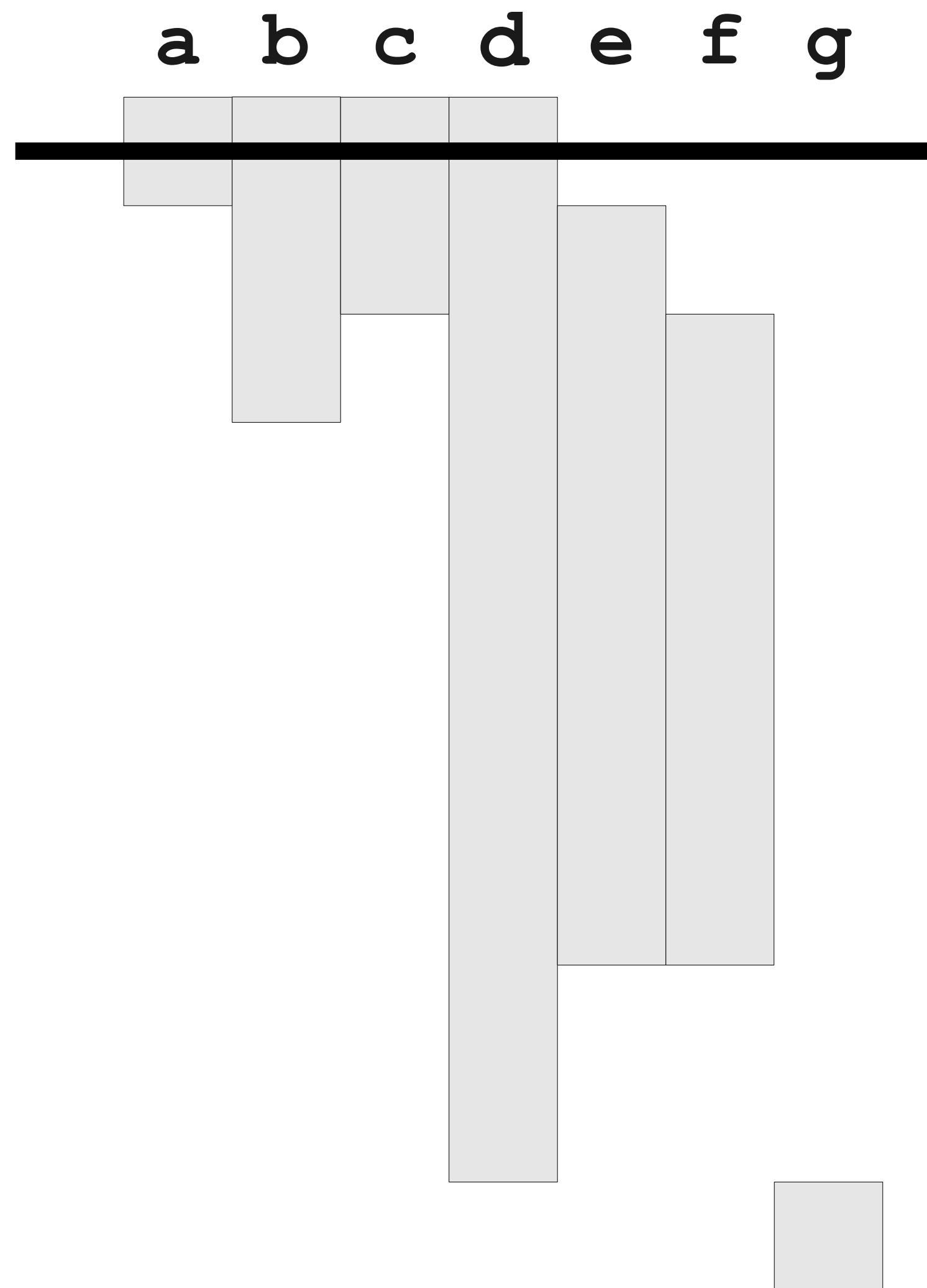
Linear-Scan Register Allocation



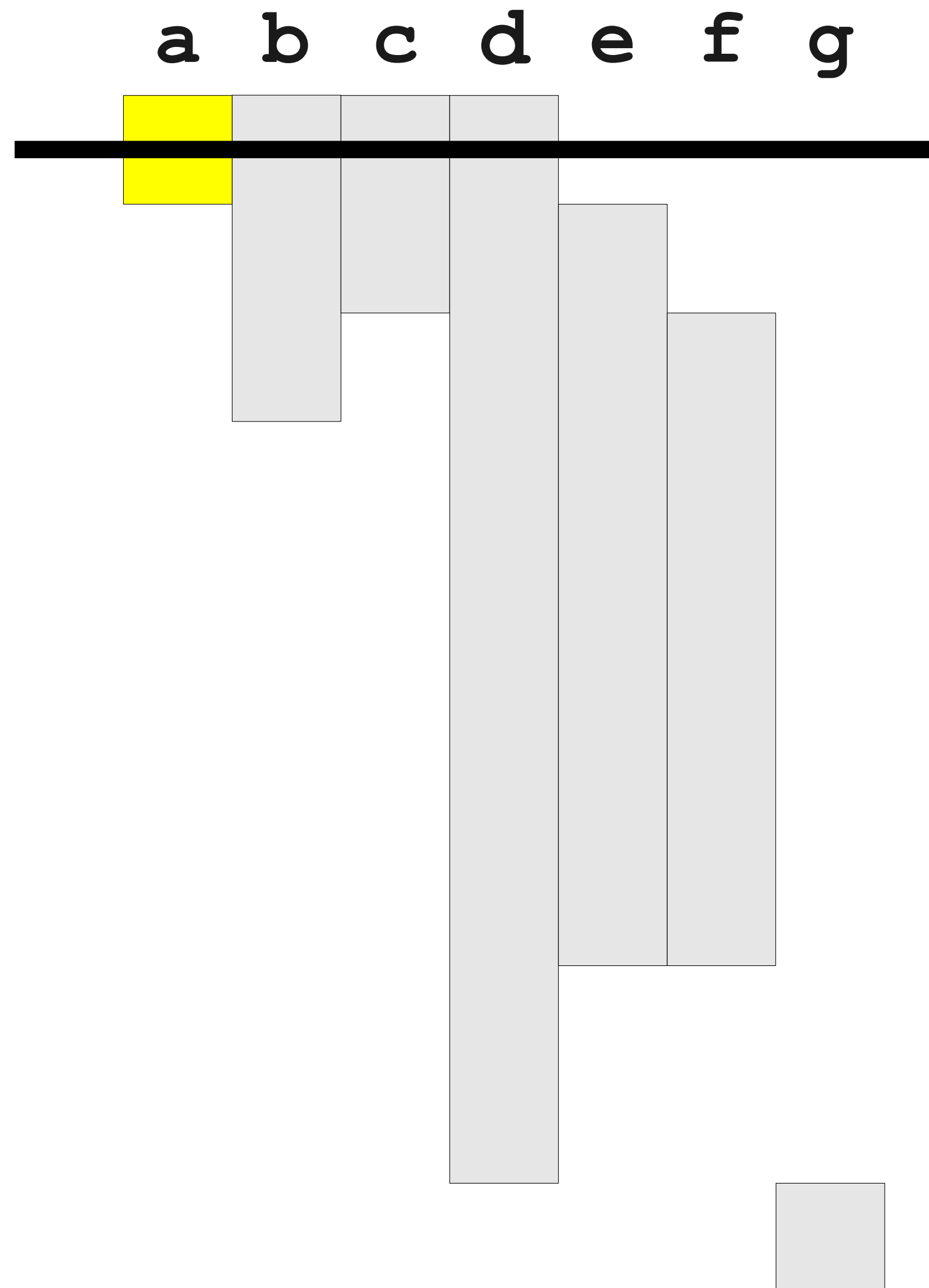
Free Registers



Linear-Scan Register Allocation



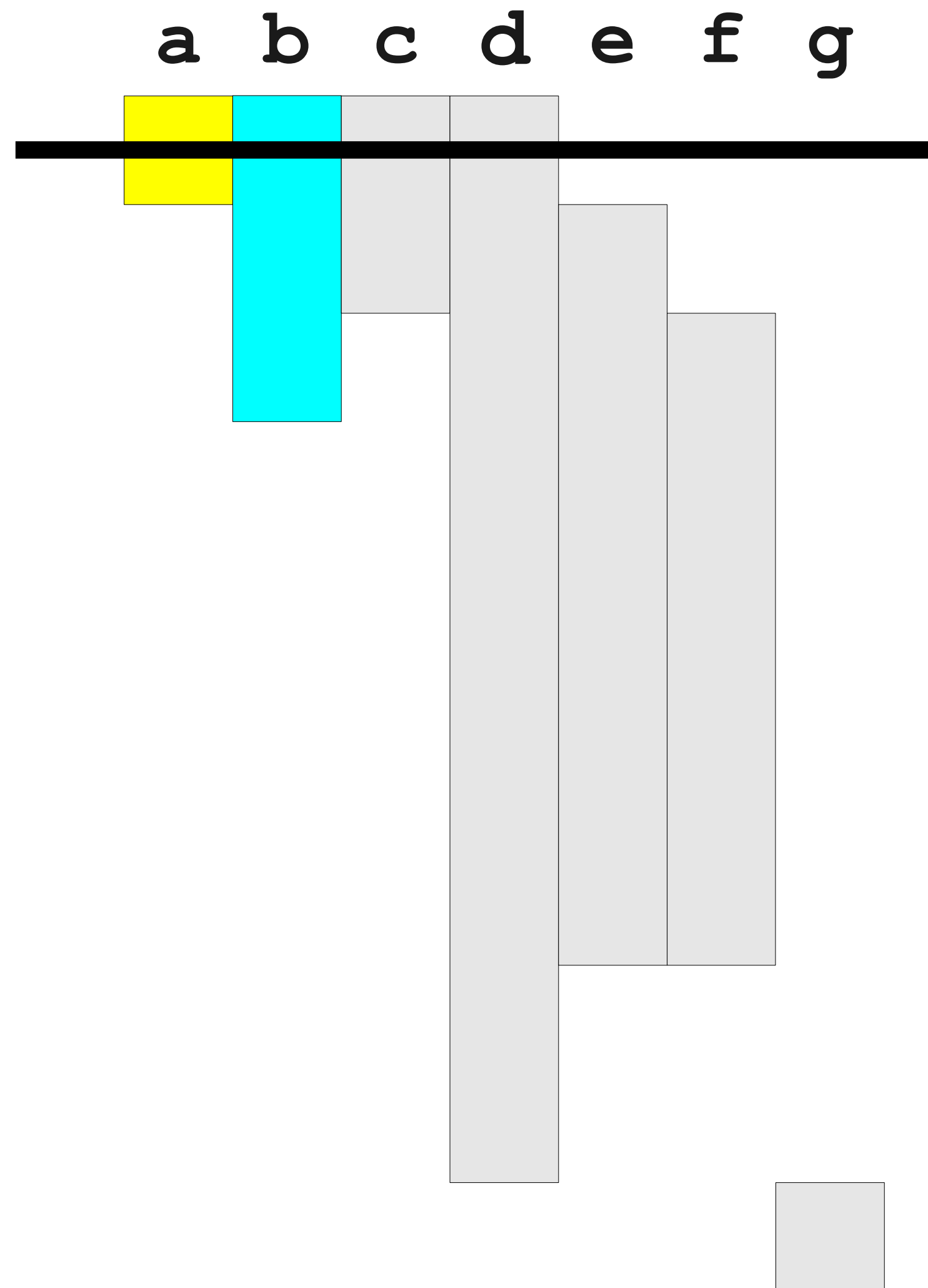
Linear-Scan Register Allocation



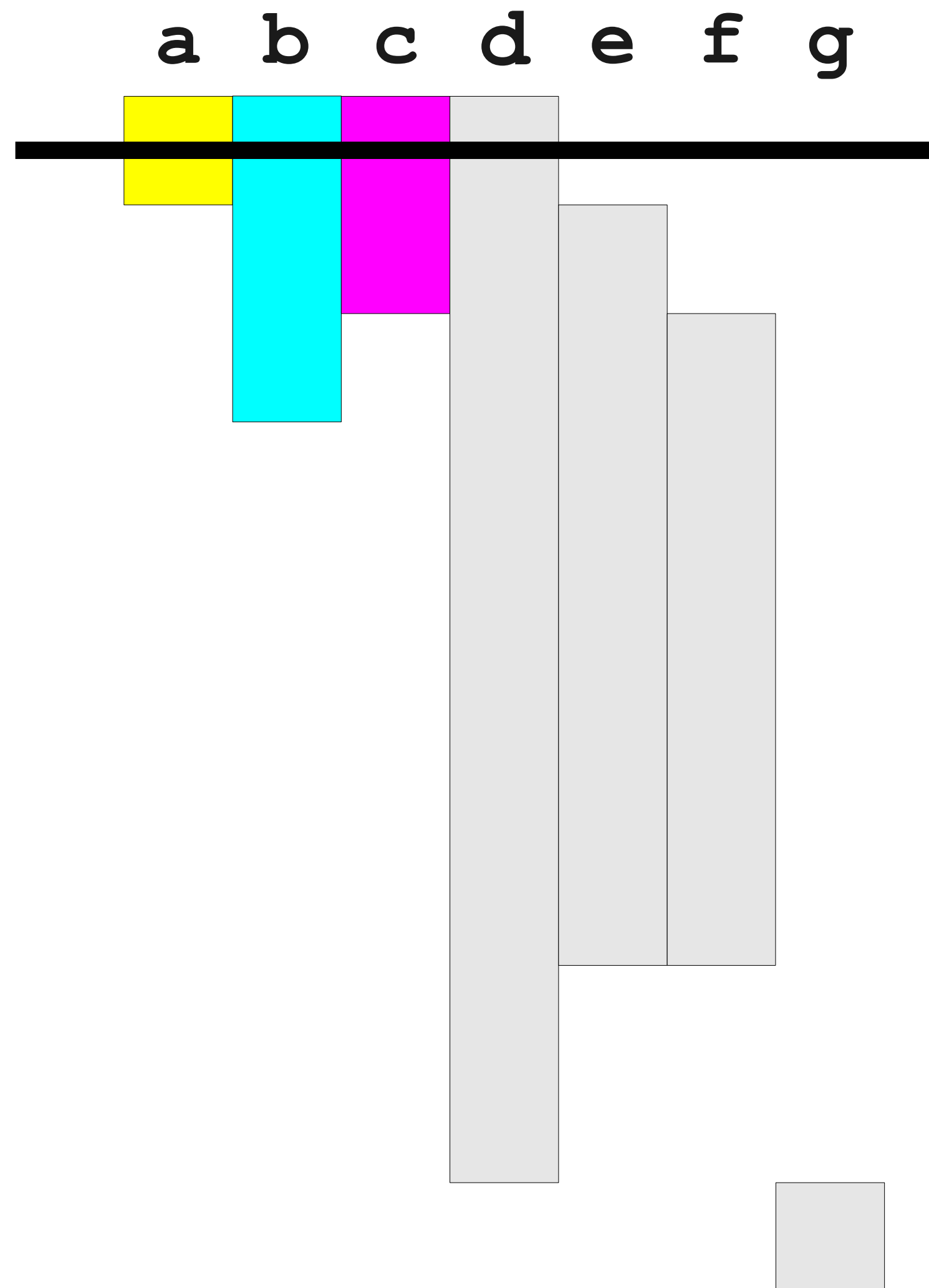
Free Registers



Linear-Scan Register Allocation



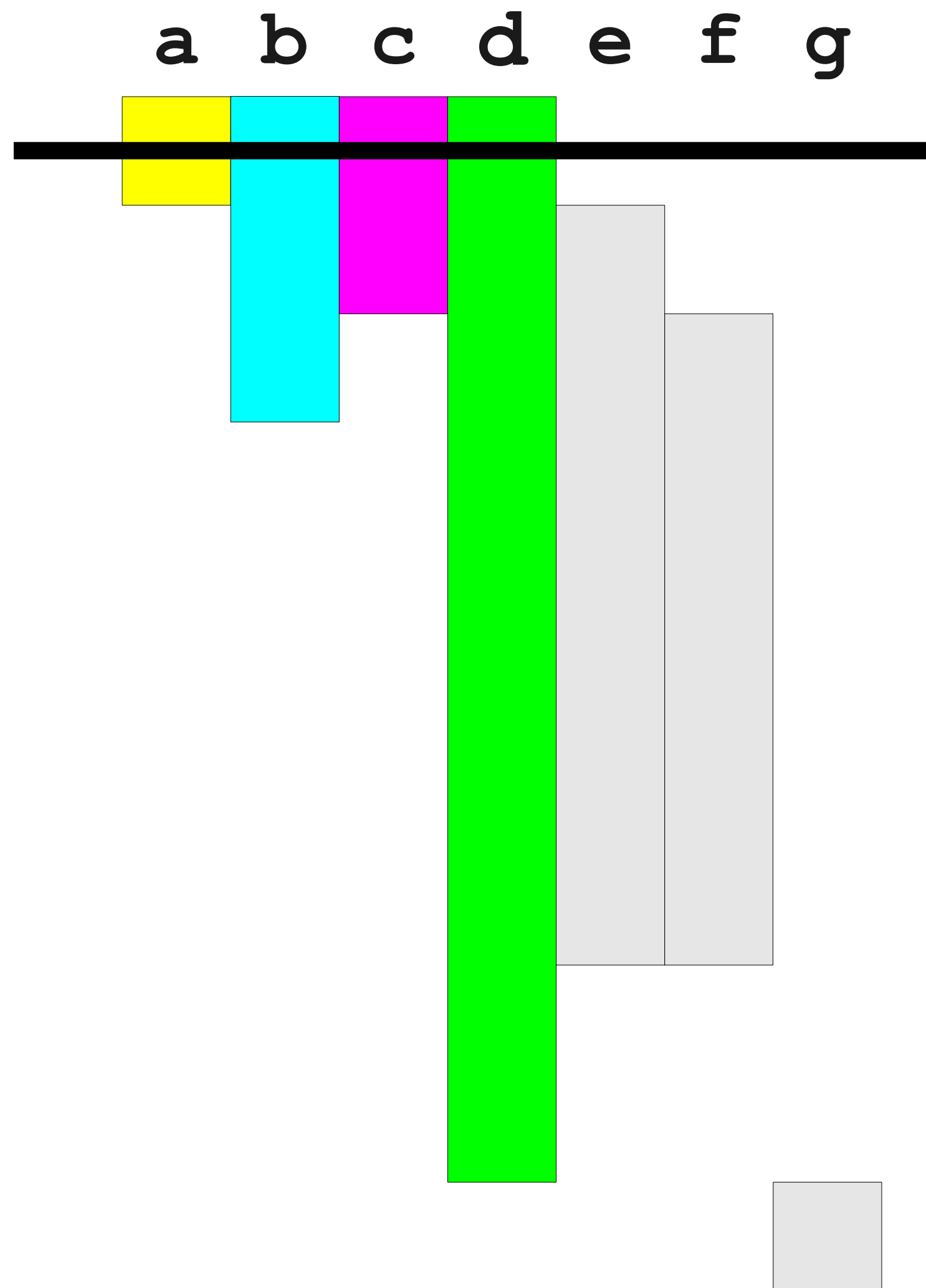
Linear-Scan Register Allocation



Free Registers

R_0	R_1	R_2	R_3
-------	-------	-------	-------

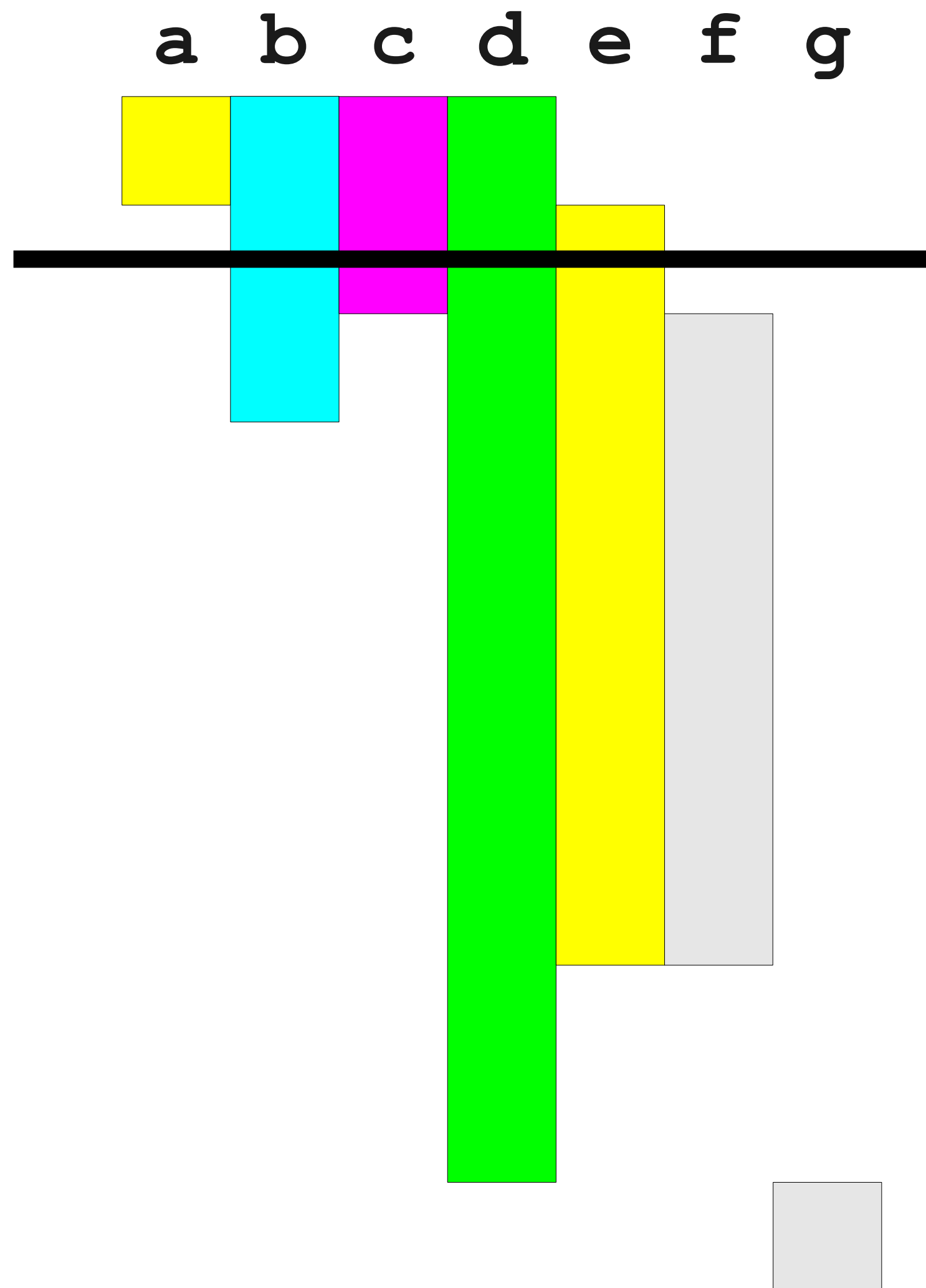
Linear-Scan Register Allocation



Free Registers

R_0	R_1	R_2	R_2
-------	-------	-------	-------

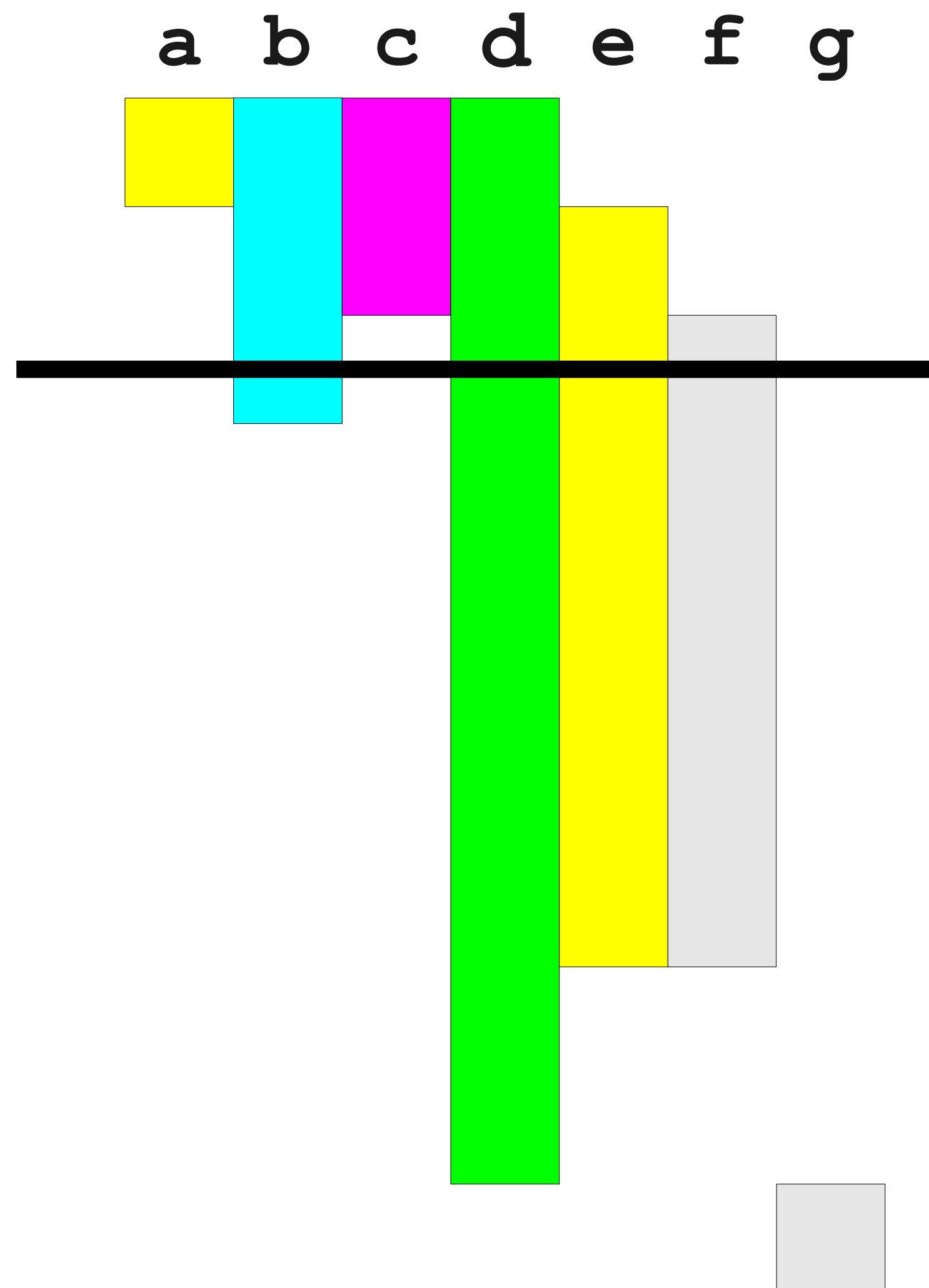
Linear-Scan Register Allocation



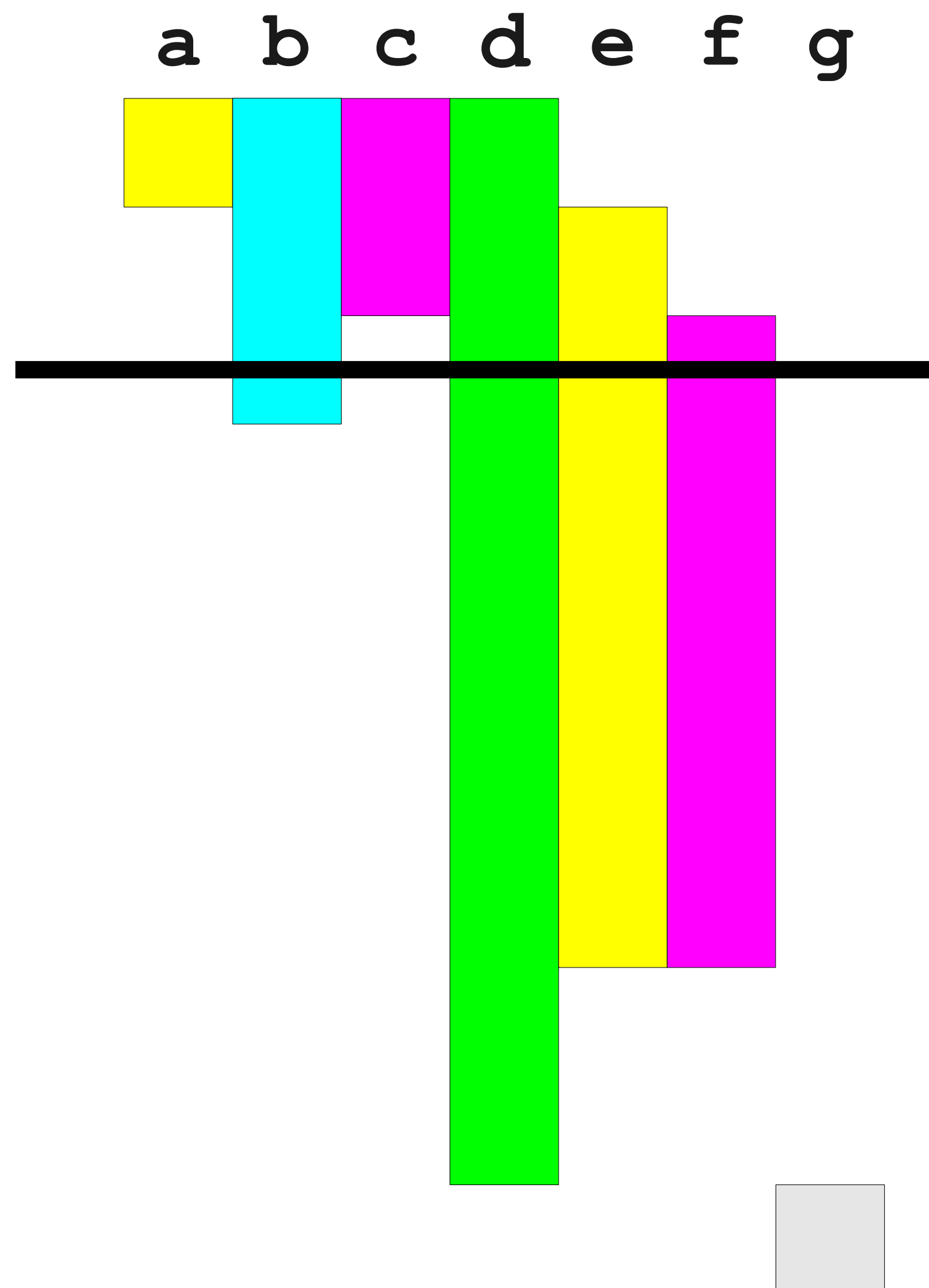
Free Registers

R_0	R_1	R_2	R_2
-------	-------	-------	-------

Linear-Scan Register Allocation

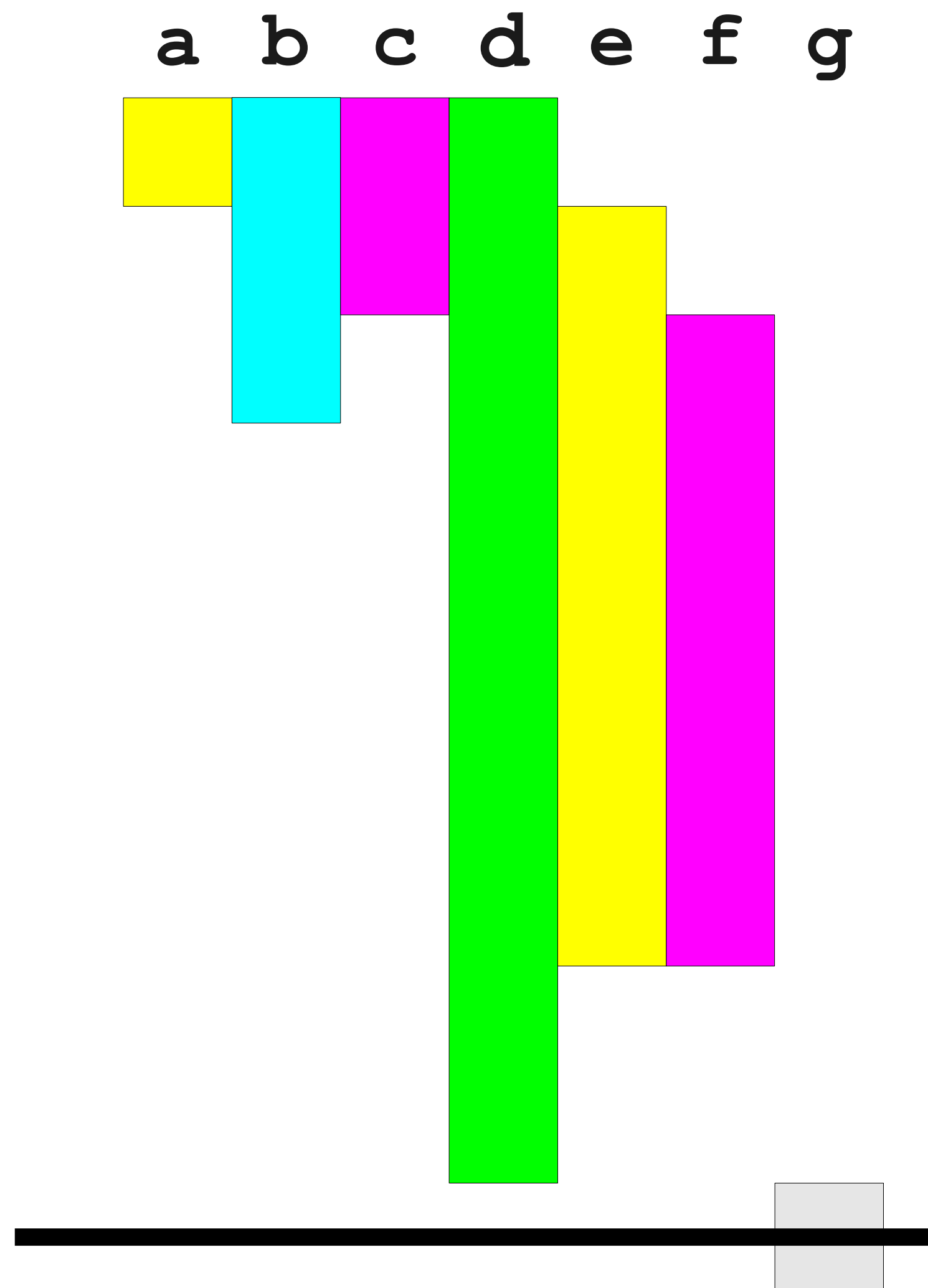


Linear-Scan Register Allocation

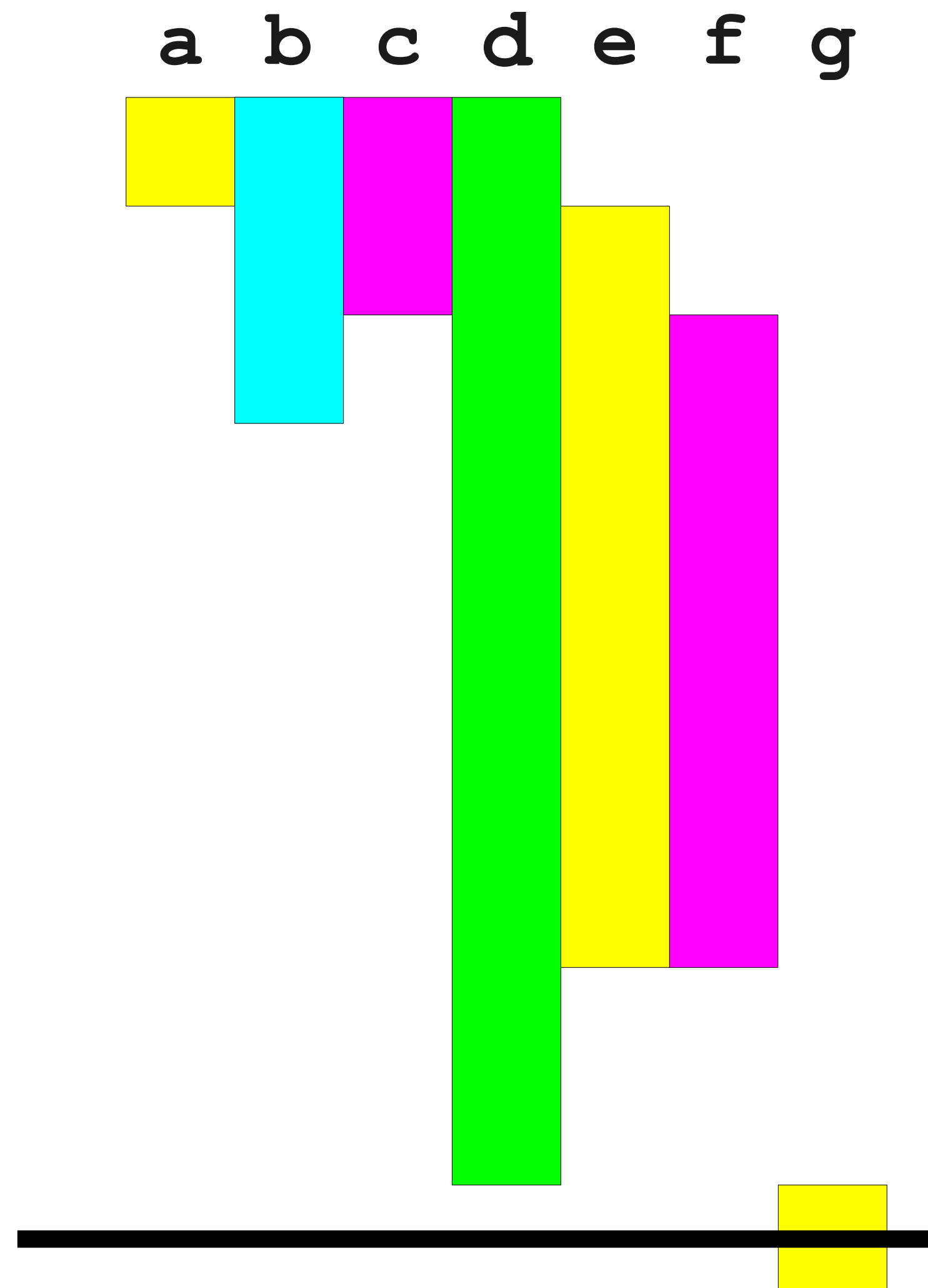


Free Registers			
R_0	R_1	R_2	R_2

Linear-Scan Register Allocation



Linear-Scan Register Allocation



Linear-Scan Register Allocation

Simple, *greedy* register-allocation strategy:

1. Compute liveness information: `live(x)`
 - recall: `live(x)` is the set of uids that are live on entry to `x`'s definition
2. Let `regs` be the set of usable registers
 - usually reserve a couple for spill code (offloading to stack) [our implementation uses `rax,rcx`]
3. Maintain “layout” `uid_loc` that maps uids to locations
 - locations include registers and stack slots `n`, starting at `n=0`
4. Scan through the program. For each instruction that defines a uid `x`
 - `used = {r | reg r = uid_loc(y) s.t. y ∈ live(x)}`
 - `available = regs - used`
 - If `available` is empty: *// no registers available, spill*
 `uid_loc(x) := slot n ; n = n + 1`
 - Otherwise, pick `r` in `available`: *// choose an available register*
 `uid_loc(x) := reg r`

Linear-Scan Register Allocation

- Advantages
 - Very efficient (after computing live intervals, runs in linear time)
 - Produces good code in many instances.
 - Allocation step works in one pass; can generate code during iteration.
 - Often used in JIT compilers like Java HotSpot.
- Pitfalls
 - Doesn't always choose a very good strategy due to *greediness*.
 - Doesn't work well with branching.

Linear Scan and Branching

```
a = d = c = 0
b = 1
if (?) {
    a = b + 1
    c = a + b
    print(a)
} else {
    d = b + 2
    c = d + b
    print(d)
}
return c
```

a = d = c = 0

// live = { }

b = 1

// live = { b }

// live = { b }

a = b + 1

d = b + 2

// live = { a, b }

// live = { d, b }

c = a + b

c = d + b

// live = { a, c }

// live = { d, c }

print(a)

print(d)

// live = { c }

// live = { c }

return c

Graph Colouring

Register Allocation via Colour Graphs

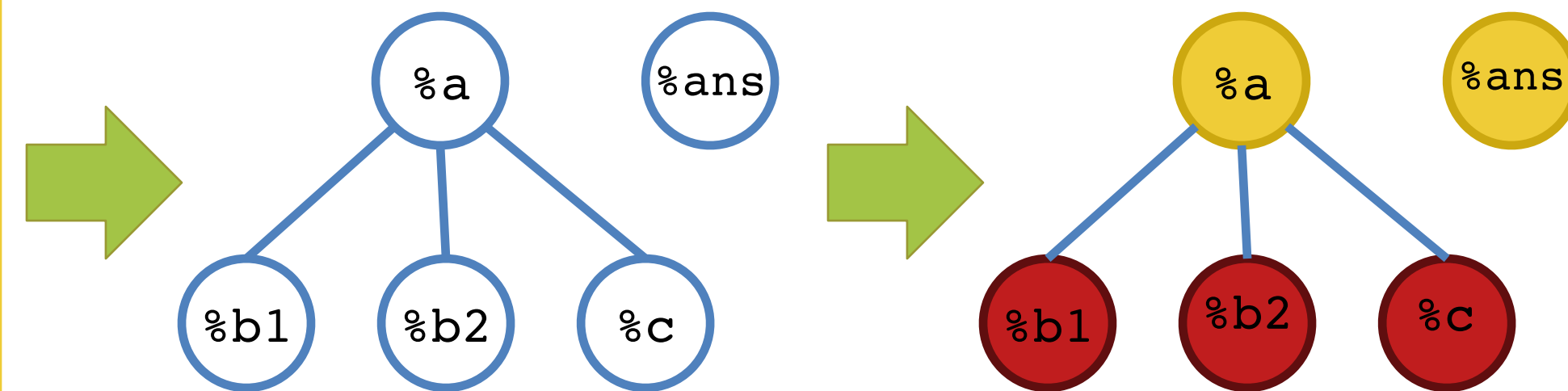
Basic process:

1. Compute liveness information for each temporary (%uid).
2. Create an *interference graph*:
 - Nodes are temporary variables (%uids).
 - There is an *edge* between node **n** and **m** if **n** is *live* at the same time as **m**
3. Try to colour the graph
 - Each colour corresponds to a register
4. In case **Step 3** fails, “spill” a temporary to the stack and repeat the whole process.
5. Rewrite the program to use registers

Interference Graphs

- Nodes of the graph are `%uids`
- Edges connect variables that *interfere* with each other
 - Two variables *interfere* if their live ranges intersect (i.e. there is an edge in the control-flow graph across which they are both live).
- Register assignment is a *graph colouring*.
 - A graph colouring assigns each node in the graph a colour (register)
 - Any two nodes connected by an edge must have different colours.
- Example:

```
// live = {%a}
%b1 = add i32 %a, 2
// live = {%a,%b1}
%c = mult i32 %b1, %b1
// live = {%a,%c}
%b2 = add i32 %c, 1
// live = {%a,%b2}
%ans = mult i32 %b2, %a
// live = {%ans}
return %ans;
```



Interference Graph

2-Colouring of the graph:
red = r8
yellow = r9

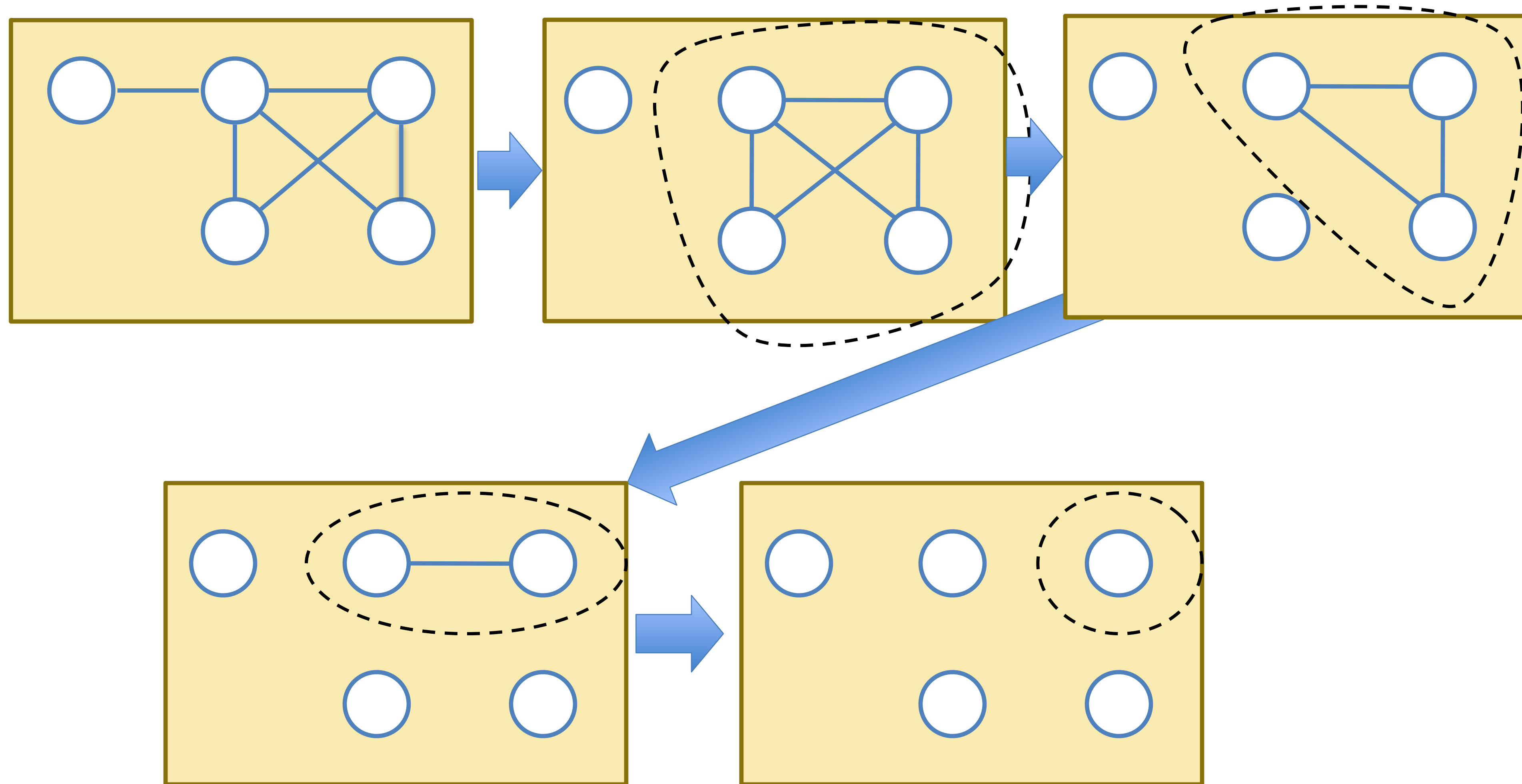
Register Allocation Questions

- Can we efficiently find a *k-colouring* of the graph whenever possible?
 - Answer: in general the problem is NP-complete (it requires search)
 - But, we can do an efficient approximation using *heuristics*.
- How do we assign registers to colours?
 - If we do this in a smart way, we can eliminate many redundant MOV instructions.
- What do we do when there aren't enough colours/registers?
 - We have to use stack space, but how do we do this effectively?

Colouring a Graph: Kempe's Algorithm

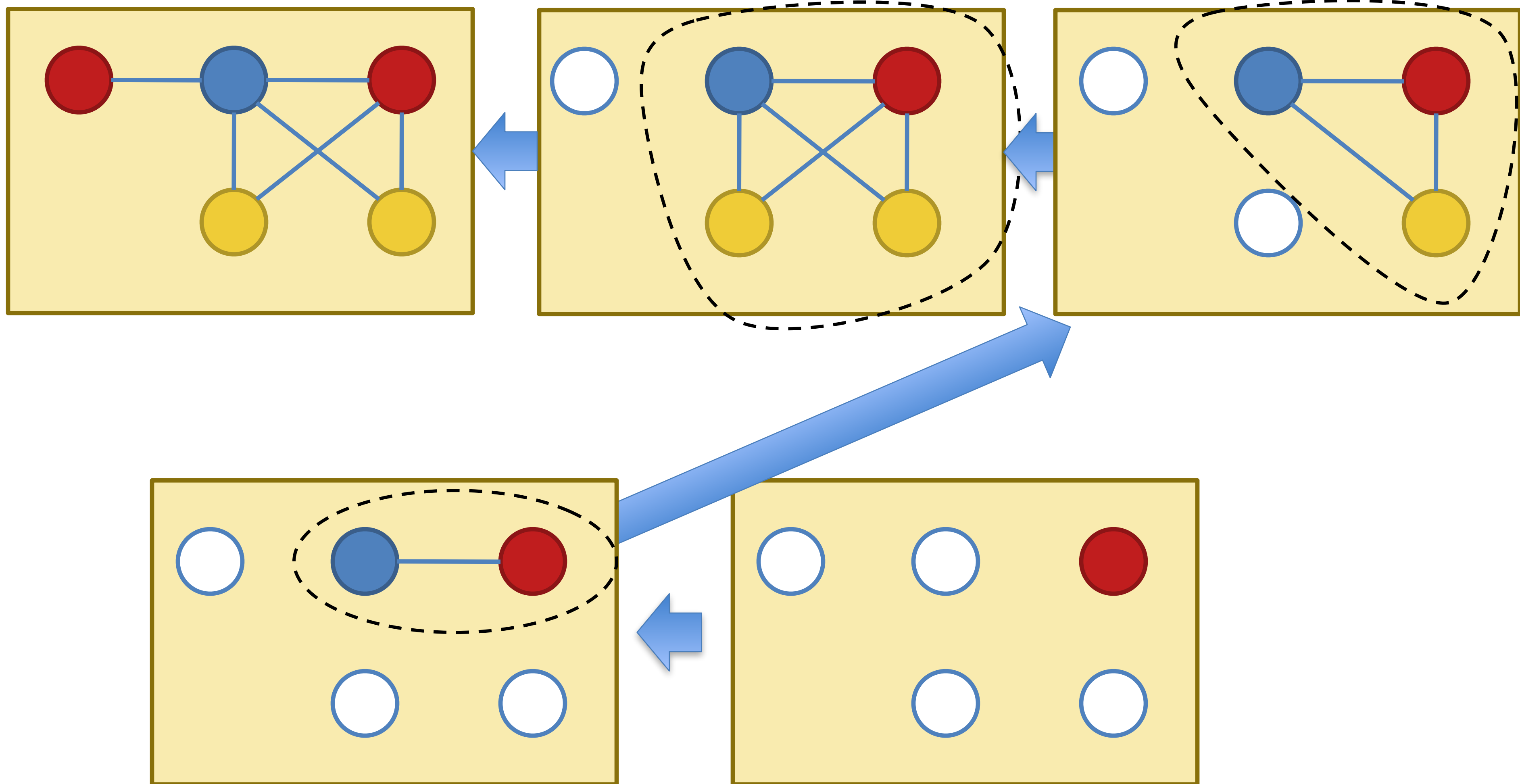
- Kempe [1879] provides this algorithm for K-coloring a graph.
- It's a recursive algorithm that works in three steps:
- Step 1: Find a node with degree $< K$ and cut it out of the graph.
 - Remove the nodes and edges.
 - This is called *simplifying* the graph
- Step 2: Recursively K-colour the remaining subgraph
- Step 3: When remaining graph is coloured, there must be at least one free colour available for the deleted node (since its degree was $< K$). Pick such a colour.

Example: 3-colour this Graph



Recurring Down the Simplified Graphs

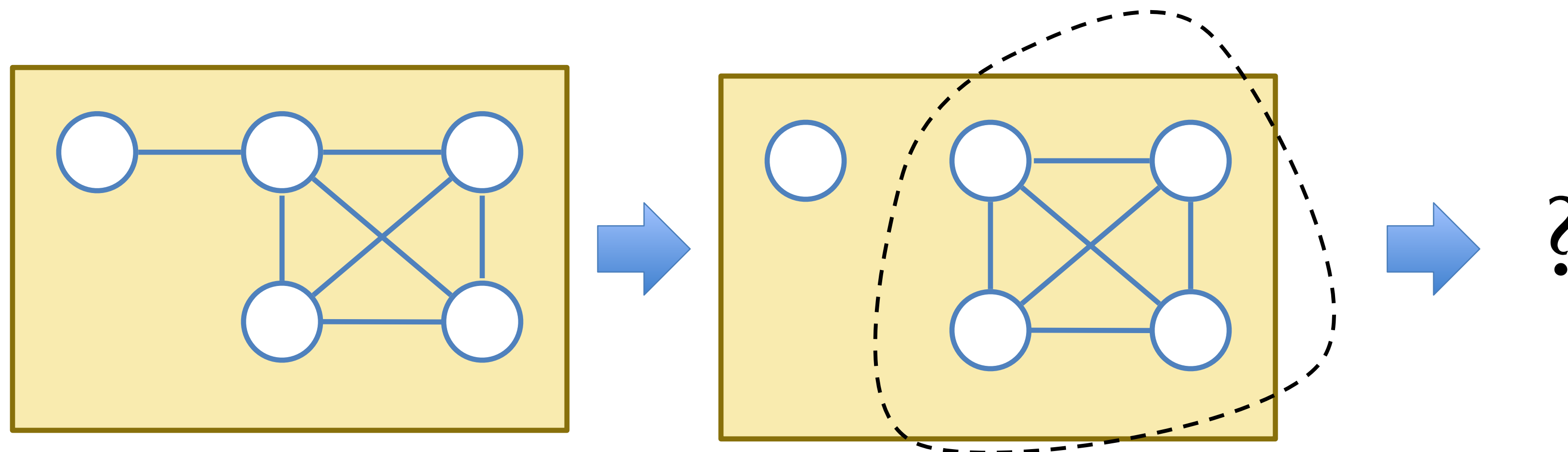
Example: 3-colour this Graph



Assigning colours on the way back up.

Failure of the Algorithm

- If the graph cannot be coloured, it will simplify to a graph where every node has at least K neighbours.
 - This can happen even when the graph is K -colourable!
 - This is a symptom of NP-hardness (it requires search)
- Example: When trying to 3-colour this graph:

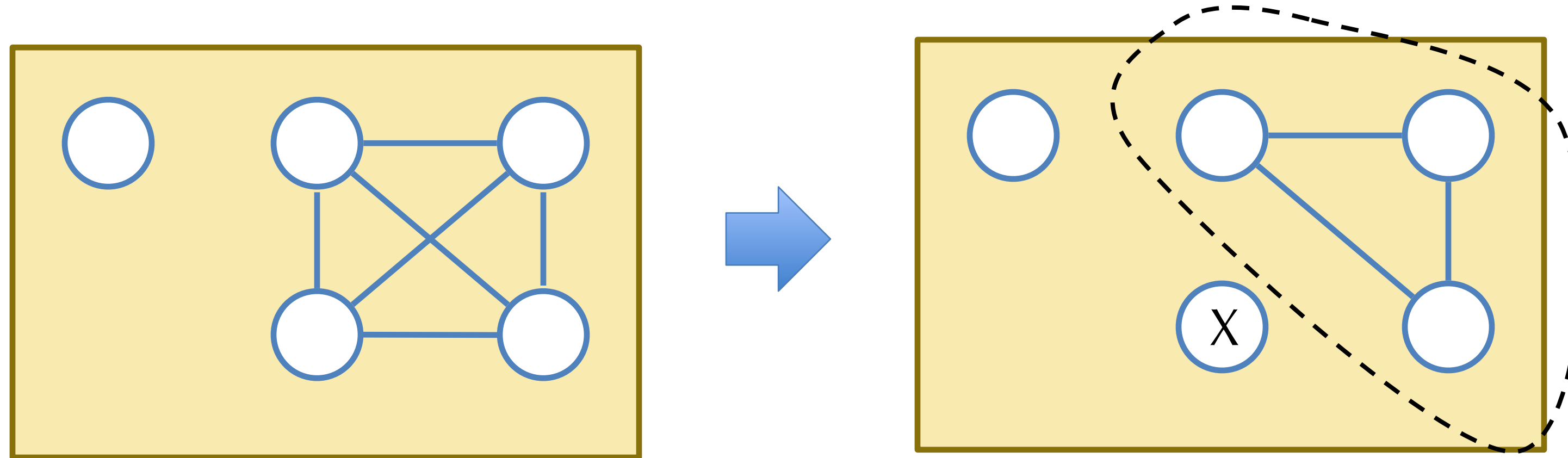


Spilling

- Idea: If we can't K-colour the graph, we need to store one temporary variable on the stack.
- Which variable to spill?
 - Pick one that isn't used very frequently
 - Pick one that isn't used in a (deeply nested) loop
 - Pick one that has high interference (since removing it will make the graph easier to colour)
- In practice: some weighted combination of these criteria
- When colouring:
 - Mark the node as spilled
 - Remove it from the graph
 - Keep recursively colouring

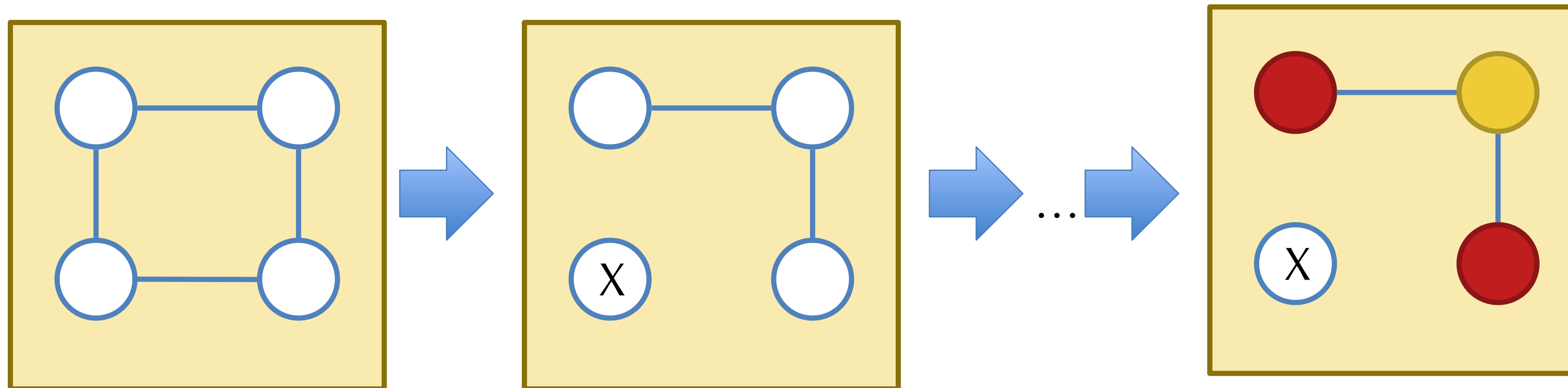
Spilling, Pictorially

- Select a node to spill
- Mark it and remove it from the graph
- Continue colouring



Optimistic Colouring

- Sometimes it is possible to colour a node marked for spilling.
 - If we get “lucky” with the choices of colours made earlier.
- Example: When 2-colouring this graph, we don't have a node with degree < 2



- Even though the node was marked for spilling, we can colour it.
- So: on the way down, mark for spilling, but don't actually spill...

Pre-coloured Nodes

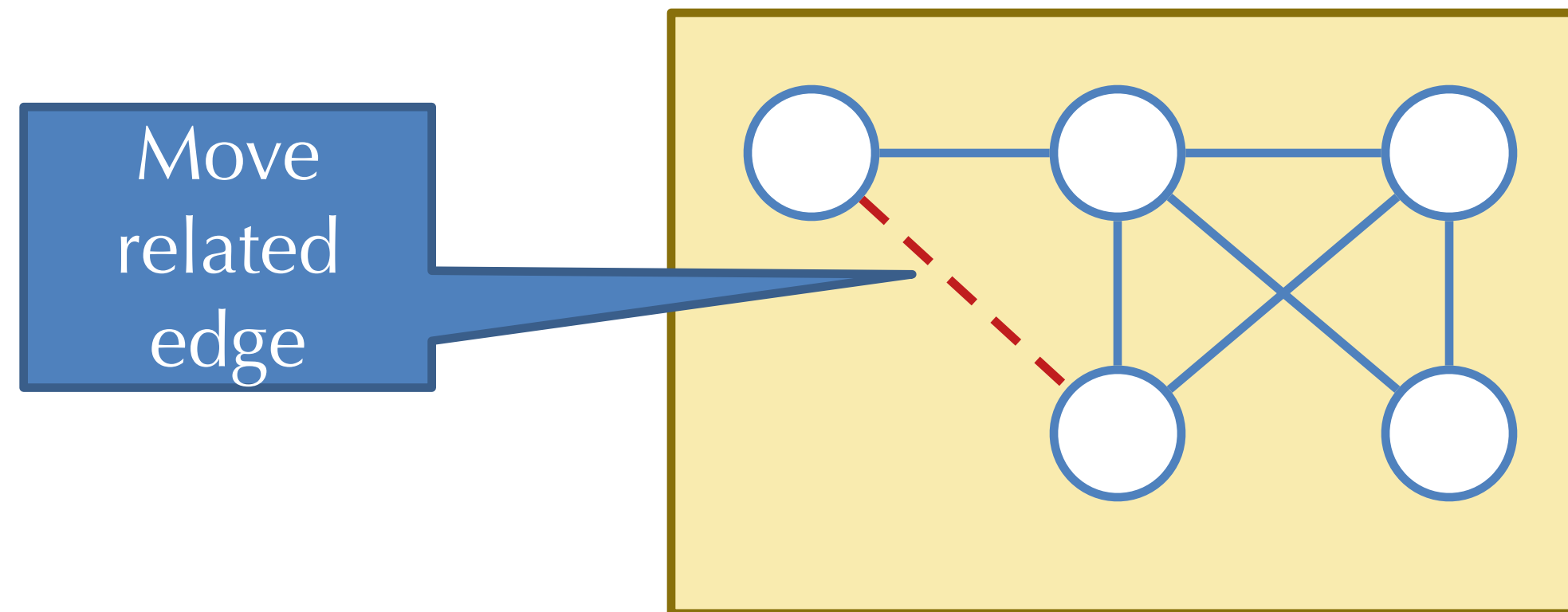
- Some variables must be pre-assigned to registers.
 - E.g. on X86 the multiplication instruction: IMul must define %rax
 - Arguments of a call have pre-assigned registers
- To properly allocate temporaries, we treat registers as nodes in the interference graph with pre-assigned colours.
 - Pre-coloured nodes can't be removed during simplification.
 - Trick: Treat pre-coloured nodes as having “infinite” degree in the interference graph – this guarantees they won't be simplified.
 - When the graph is empty except the pre-coloured nodes, we have reached the point where we start colouring the rest of the nodes.

Picking Good Colours

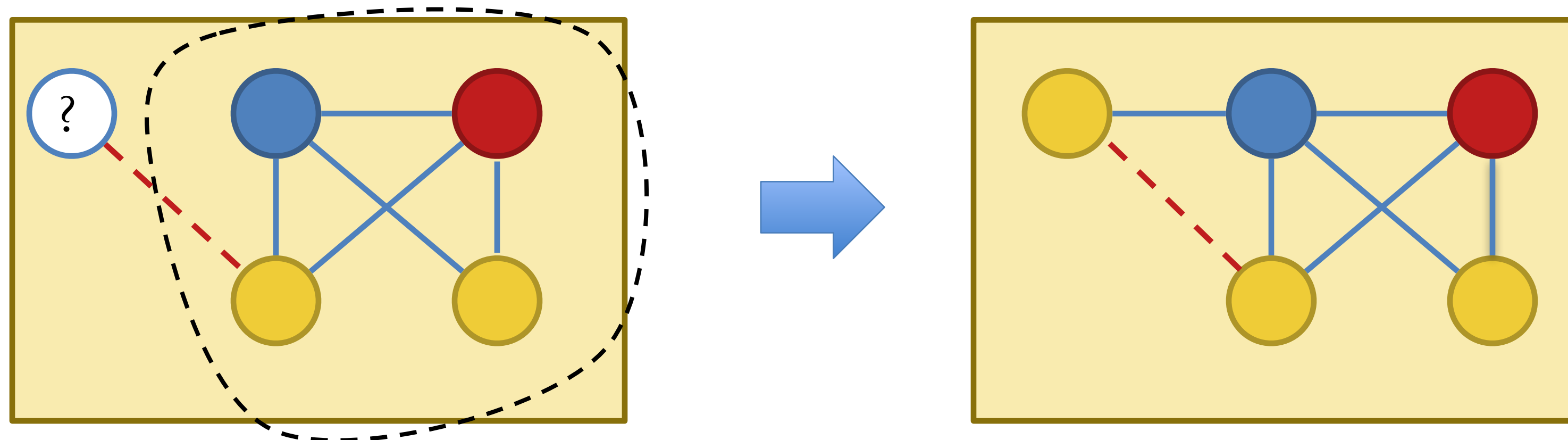
- When choosing colours during the colouring phase, *any* choice is semantically correct, but some choices are better for performance.
- Example:
 `%t1 = %t2`
 - If t1 and t2 can be assigned the same register (colour) then this move is redundant and can be eliminated.
- A simple colour choosing strategy that helps eliminate such moves:
 - Add a new kind of “move related” edge between the nodes for t1 and t2 in the interference graph.
 - When choosing a colour for t1 (or t2), if possible pick a colour of an already coloured node reachable by a move-related edge.

Example Colour Choice

- Consider 3-colouring this graph, where the dashed edge indicates that there is a Move from one temporary to another.

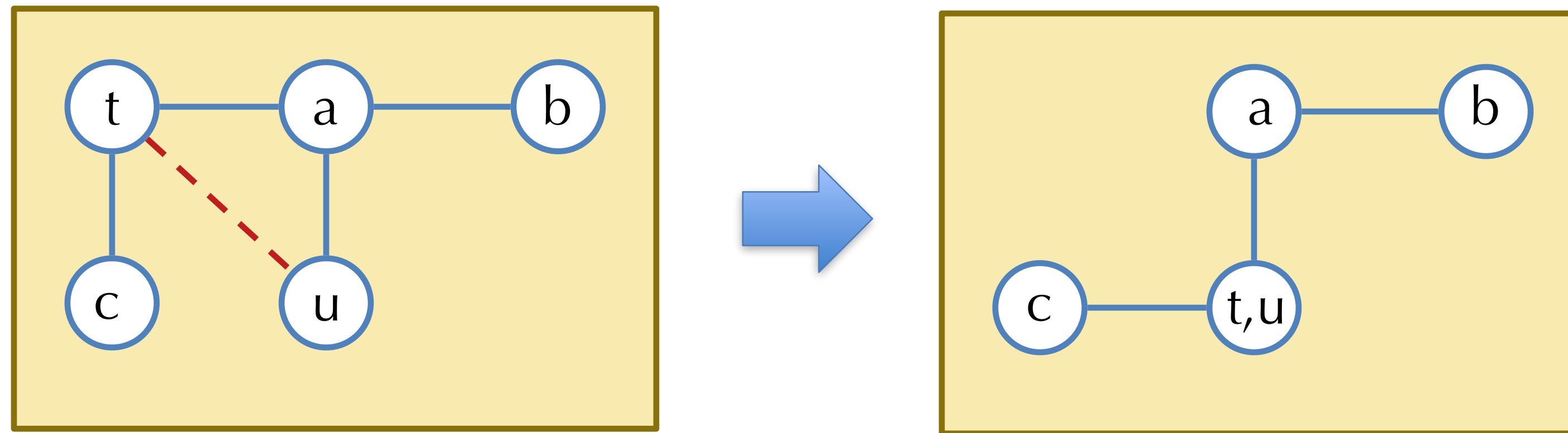


- After colouring the rest, we have a choice:
 - Picking yellow is better than red because it will eliminate a move.

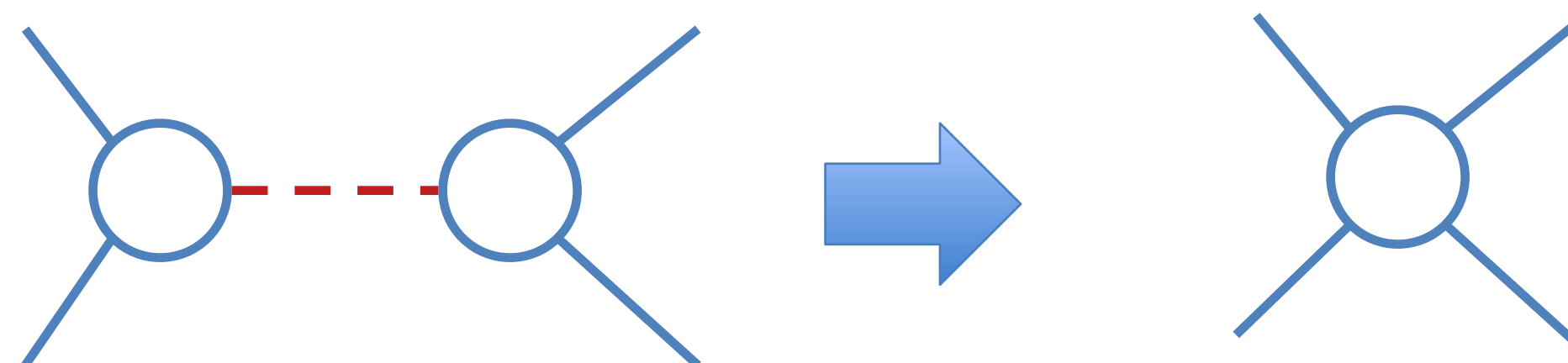


Coalescing Interference Graphs

- A more aggressive strategy is to *coalesce* nodes of the interference graph if they are connected by move-related edges.
 - Coalescing the nodes *forces* the two temporaries to be assigned the same register.



- Idea: interleave simplification and coalescing to maximize the number of moves that can be eliminated.
- Problem: coalescing can sometimes increase the degree of a node.



Conservative Coalescing

- Two strategies are guaranteed to preserve the k -colorability of the interference graph.
- *Brigg's strategy*: It's safe to coalesce \mathbf{x} & \mathbf{y} if the resulting node will have fewer than k neighbours (with degree $\geq k$).
- *George's strategy*: We can safely coalesce \mathbf{x} & \mathbf{y} if for every neighbour \mathbf{t} of \mathbf{x} , either \mathbf{t} already interferes with \mathbf{y} or \mathbf{t} has degree $< k$.

Graph Colouring Algorithm

1. Build interference graph (pre-colour nodes as necessary).
 - Add move related edges
2. Reduce the graph (building a stack of nodes to color).
 - a. Simplify the graph as much as possible without removing nodes that are move-related (i.e. have a move-related neighbour). Remaining nodes are high degree or move-related.
 - b. Coalesce move-related nodes using Brigg's or George's strategy.
 - c. Coalescing can reveal more nodes that can be simplified, so repeat 2.a and 2.b until no node can be simplified or coalesced.
 - d. If no nodes can be coalesced *freeze* (remove) a move-related edge and keep trying to simplify/coalesce.
3. If there are non-precoloured nodes left, mark one for spilling, remove it from the graph and continue doing step 2.
4. When only pre-coloured node remain, start colouring (popping simplified nodes off the top of the stack).
 - a. If a node must be spilled, insert spill code as on slide “*Example Spill Code*” and rerun the whole register allocation algorithm, starting at step 1.

Was it worth it?

Demo: Register allocation in HW6

For HW6

- HW 6 implements two naive register allocation strategies:
- `no_reg_layout`: spill all registers to the stack
- `greedy_layout`: puts the first few uids in available registers and spills the rest. It uses liveness information to recycle available registers when their current value becomes dead (see the slides above).
- Your job: do “better” than these via graph colouring.
- Quality Metric:
 - the total number of memory accesses in x86 program, which is the sum of:
 - the number of Ind2 and Ind3 operands
 - the number of Push and Pop instructions
 - shorter code is better

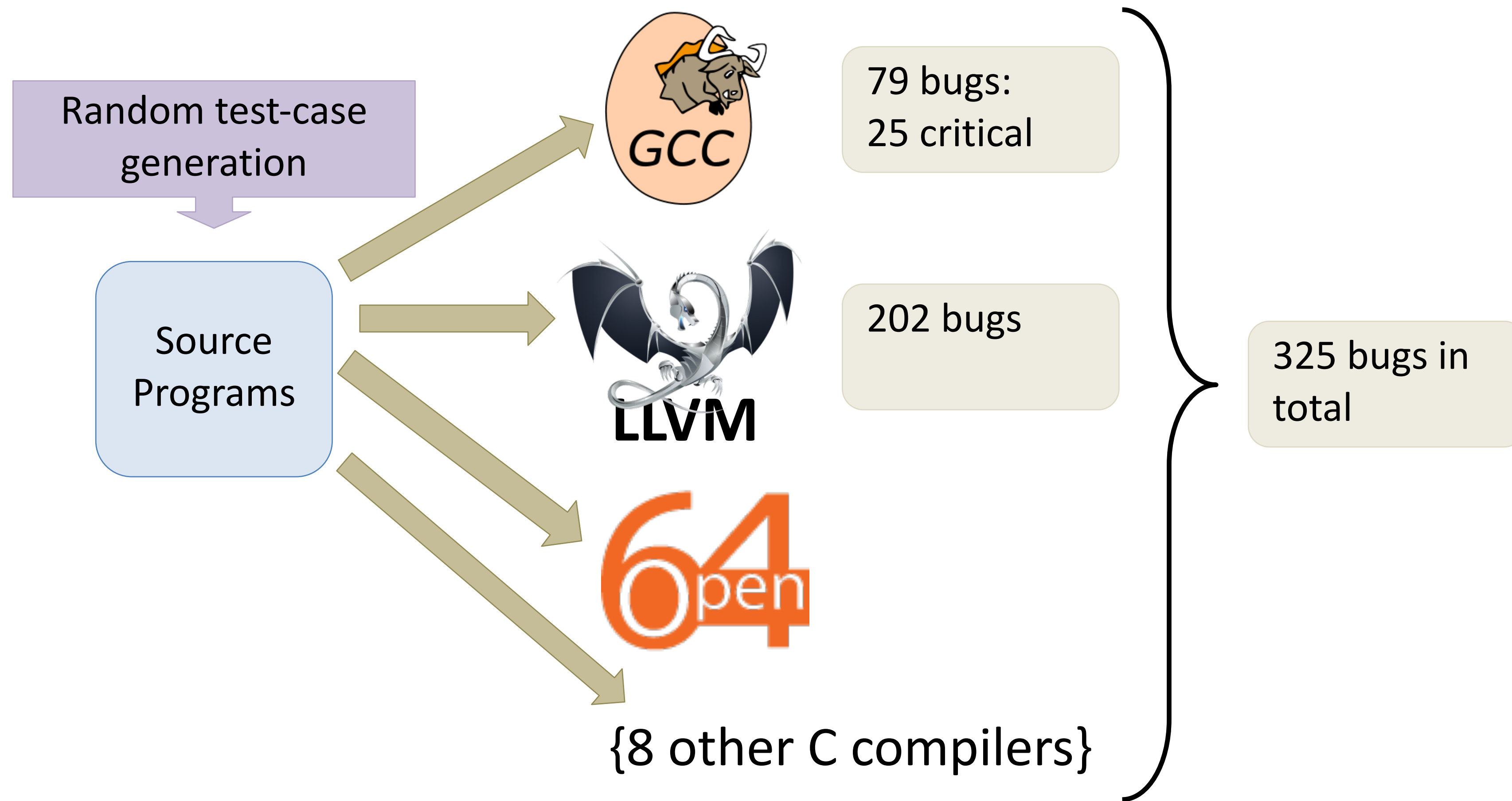
Current Research in Programming Languages and Compilers

Validating Compilers

- The job of a compiler is to *translate from the syntax of one language to another, but preserve the semantics*.
- Compiler correctness is critical
 - Trustworthiness of every component built in a compiled language depends on trustworthiness of the compiler
- Compilers tend to be well-engineered and well-tested, but that does not mean they are *bug-free*.

Testing Compilers

Finding and Understanding Bugs in C Compilers. Yang et al. PLDI 2011



Finding and Understanding Bugs in C Compilers

Xuejun Yang Yang Chen Eric Eide John Regehr

University of Utah, School of Computing
{jxyang, chenyang, eeide, regehr}@cs.utah.edu

(in PLDI 2011)

Compilers should be correct.

To improve the quality of C compilers, we created Csmith, a **randomized test-case generation tool**, and spent **three years** using it to find compiler bugs.

During this period we reported **more than 325 previously unknown bugs** to compiler developers.

The striking thing about our **CompCert** results is that the middle-end bugs we found in all other compilers are **absent**.

As of early 2011, the under-development version of **CompCert** is the only compiler we have tested for which Csmith cannot find wrong-code errors. This is not for lack of trying: we have devoted about six CPU-years to the task.

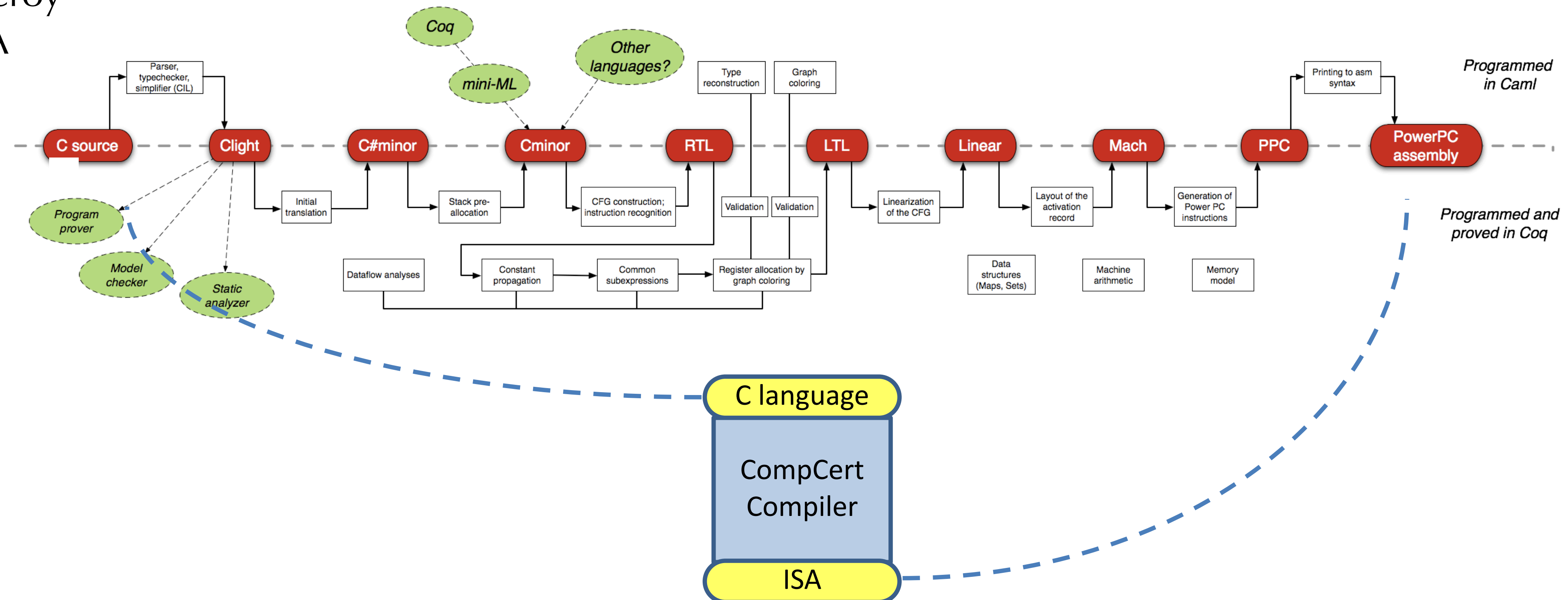
The apparent unbreakability of CompCert supports a strong argument that developing compiler optimizations within a proof framework, where safety checks are explicit and machine-checked, has tangible benefits for compiler users.

Verified Compilation



Xavier Leroy
INRIA

CompCert (2006-now):
Optimising C Compiler,
proved correct *end-to-end*
with machine-checked proof in Coq



Comparing Behaviours

- Consider two programs P1 and P2 possibly in different languages.
 - e.g. P1 is an Oat program, P2 is its compilation to LL
- The semantics of the languages associate to each program a set of observable behaviours:

$\mathbf{B}(P1)$ and $\mathbf{B}(P2)$

- Note: $|\mathbf{B}(P)| = 1$ if P is deterministic, > 1 otherwise

What is Observable?

- For C-like languages:

observable behavior ::=

- | terminates(st) (i.e. observe the final state)
- | diverges
- | goeswrong

- For pure functional languages:

observable behavior ::=

- | terminates(v) (i.e. observe the final value)
- | diverges
- | goeswrong

What about I/O?

- Add a *trace* of input-output events performed:

t	$::= [] \mid e :: t$	(finite traces)
coind. T	$::= [] \mid e :: T$	(finite and infinite traces)

observable behavior $::=$

terminates(t, st)	(end in state st after trace t)
diverges(T)	(loop, producing trace T)
goeswrong(t)	

Examples

- P1:
print(1); / st \Rightarrow terminates(out(1)::[],st)
- P2:
print(1); print(2); / st \Rightarrow terminates(out(1)::out(2)::[],st)
- P3:
WHILE true DO print(1) END / st \Rightarrow diverges(out(1)::out(1)::....)
- So **B(P1) \neq B(P2) \neq B(P3)**

Bisimulation

- Two programs P1 and P2 are *bisimilar* whenever:

$$\mathbf{B}(P1) = \mathbf{B}(P2)$$

- The two programs are completely indistinguishable.
- But... this is often too strong in practice.

Compilation Reduces Nondeterminism

- Some languages (like C) have underspecified behaviours:
 - Example: order of evaluation of expressions $f() + g()$

- Concurrent programs often permit nondeterminism
 - Classic optimizations can reduce this nondeterminism
 - Example:

$a := x + 1; b := x + 1 \quad || \quad x := x + 1$

vs.

$a := x + 1; b := a \quad || \quad x := x + 1$

- LLVM explicitly allows nondeterminism:
 - undef values (not part of LLVM lite)

Backward Simulation

- Program P2 can exhibit fewer behaviours than P1:

$$\mathbf{B}(P1) \supseteq \mathbf{B}(P2)$$

- All of the behaviours of P2 are permitted by P1, though some of them may have been eliminated.
- Also called *refinement*.

Related Research Topics

Automated Parallelisation

- Moore's law: processor advances double speed every 18 months
- Moore's law ended in 2006 for single-threaded applications
- Started to hit fundamental limits in how small transistors can be
- Processor manufacturers shifted to multi-core processors
- Need new compiler technology to take advantage of multi-core – automatically find and exploit opportunities for parallel execution

Program Analysis

- The goal of a program analysis is to *answer questions about the run-time behaviour of software*
- In compilers: *data flow analysis, control flow analysis*
 - Typical goal: determine whether an optimisation is safe
- Research in program analysis has shifted to more sophisticated properties:
 - *Numerical analyses*, e.g., find geometric regions that contain reachable values for integer variables.
Can be used to verify absence of buffer overflows.
 - *Shape analyses* – determine whether a data structure in the heap is a list, a tree, a graph,...
Can be used to verify memory safety.
 - *Resource analyses* – e.g., find a conservative upper bound on the run-time complexity of a loop.
Can be used to find timing side-channel attacks.
 - *Concurrency analysis*: find all data races in a multi-threaded program.
- Industrial program analysis:
 - Static Driver Verifier (Microsoft): finds bugs in device driver code
 - Infer (Facebook): proves memory safety & finds race conditions
 - Astrée (AbsInt): static analyser for safety-critical embedded code (e.g., automotive & aerospace applications)

Program Verification and Synthesis

- *Verification*: Given a program and a specification, prove that the program satisfies the specification
- *Synthesis*: Given a specification, find a program that satisfies the specification
 - Kind of a “compilation on steroids” from language of specifications to a programming language

```
void swap(loc x, loc y)
```

$\{ x \mapsto a \wedge y \mapsto b \}$

```
void swap(loc x, loc y)
```

$\{ x \mapsto a \wedge y \mapsto b \}$

`void swap(loc x, loc y)`

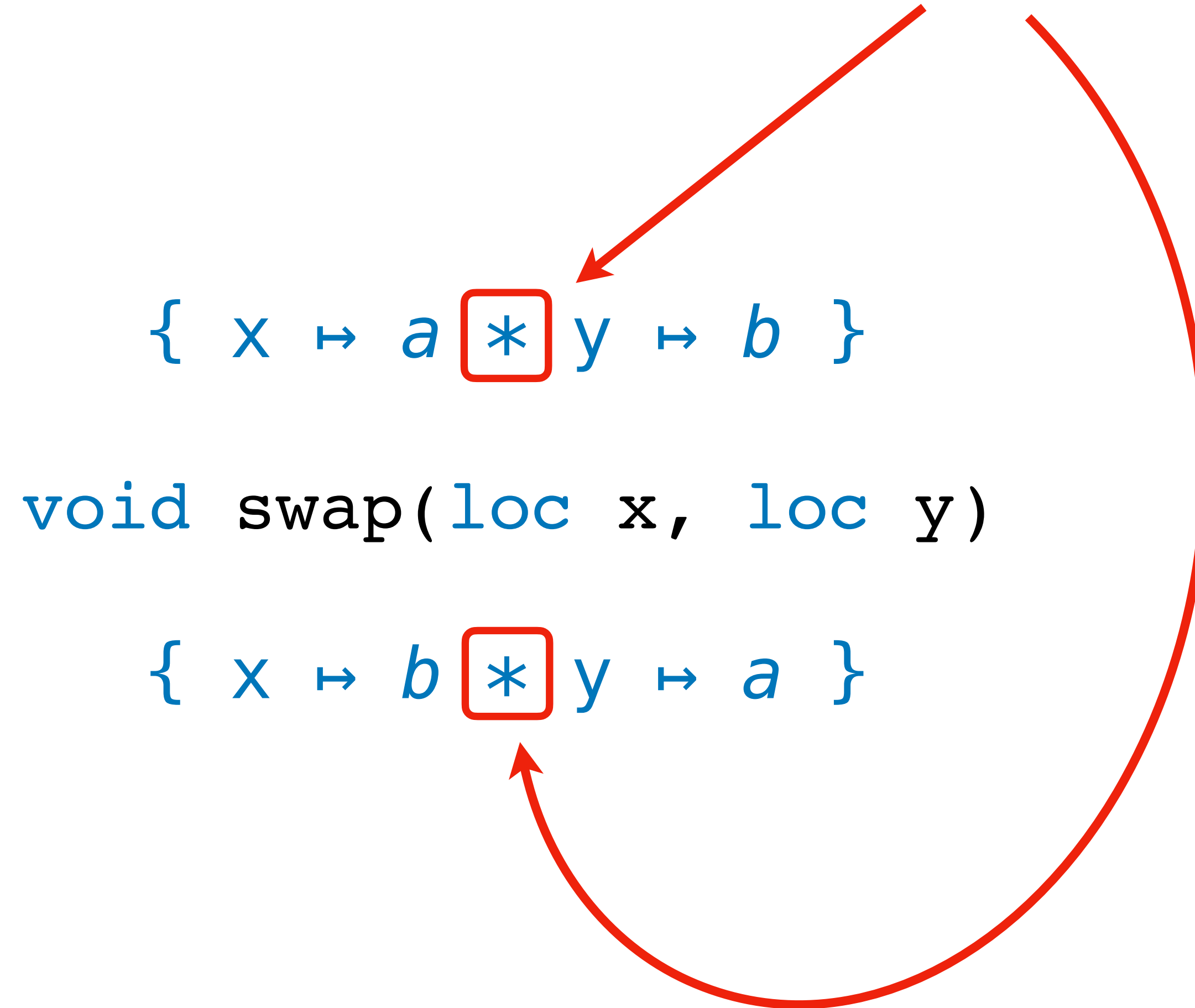
$\{ x \mapsto b \wedge y \mapsto a \}$

“x and y are different memory locations”

{ $x \mapsto a$ * $y \mapsto b$ }

void swap(loc x, loc y)

{ $x \mapsto b$ * $y \mapsto a$ }



$\{ \boxed{x} \mapsto a * \boxed{y} \mapsto b \}$

`void swap(loc \boxed{x} , loc \boxed{y})`

$\{ \boxed{x} \mapsto b * \boxed{y} \mapsto a \}$

$\{ x \mapsto \boxed{a} * y \mapsto \boxed{b} \}$

`void swap(loc x, loc y)`

$\{ x \mapsto \boxed{b} * y \mapsto \boxed{a} \}$

$$\{ x \mapsto \boxed{a} * y \mapsto b \}$$

??

$$\{ x \mapsto b * y \mapsto \boxed{a} \}$$

let a2 = *x;

{ x ↦ a2 * y ↦ b }

??

{ x ↦ b * y ↦ a2 }

```
let a2 = *x;
```

```
let b2 = *y;
```

```
{ x ↦ a2 * y ↦ b2 }
```

??

```
{ x ↦ b2 * y ↦ a2 }
```

```
let a2 = *x;
```

```
let b2 = *y;
```

```
*x = b2;
```

```
{ x ↦ b2 * y ↦ b2 }
```

??

```
{ x ↦ b2 * y ↦ a2 }
```

```
let a2 = *x;
```

```
let b2 = *y;
```

```
*x = b2;
```

```
*y = a2;
```

```
{ x ↦ b2 * y ↦ a2 }
```

??

```
{ x ↦ b2 * y ↦ a2 }
```

```
let a2 = *x;
```

```
let b2 = *y;
```

```
*x = b2;
```

```
*y = a2;
```

```
{ x ↦ b2 * y ↦ a2 }
```

??

```
{ x ↦ b2 * y ↦ a2 }
```

```
x ↦ b2 * y ↦ a2 ⊢ x ↦ b2 * y ↦ a2
```

```
let a2 = *x;
```

```
let b2 = *y;
```

```
*x = b2;
```

```
*y = a2;
```

```
{ x ↦ b2 * y ↦ a2 }
```

??

```
{ x ↦ b2 * y ↦ a2 }
```

$x \mapsto b2 * y \mapsto a2 \vdash x \mapsto b2 * y \mapsto a2$



```
void swap(loc x, loc y) {  
    let a2 = *x;  
    let b2 = *y;  
    *x = b2;  
    *y = a2;  
}
```


Transforming Entailment

$$P \rightsquigarrow Q$$

There exists a program **c**, such that
for *any* initial state satisfying **P**,
c, after it terminates,
will transform to a state satisfying **Q**.

$$x \mapsto a \rightsquigarrow x \mapsto 42$$

“Proof”: $*x = 42$

$$\Gamma; \{ \text{emp} \} \rightsquigarrow \{ \text{emp} \} \mid \mathbf{skip} \quad (\text{Emp})$$

$$\frac{\begin{array}{c} a \in \text{GV}(\Gamma, P, Q) \quad y \text{ is fresh} \\ \Gamma, y; [y/a]\{ x \mapsto y * P \} \rightsquigarrow [y/a]\{ Q \} \mid c \end{array}}{\Gamma; \{ x \mapsto a * P \} \rightsquigarrow \{ Q \} \mid \mathbf{let } y = *x; c} \quad (\text{Read})$$

$$\frac{\begin{array}{c} \text{EV}(\Gamma, P, Q) \cap \text{Vars}(R) = \emptyset \\ \Gamma; \{ P \} \rightsquigarrow \{ Q \} \mid c \end{array}}{\Gamma; \{ P * R \} \rightsquigarrow \{ Q * R \} \mid c} \quad (\text{Frame})$$

$$\frac{\begin{array}{c} \text{Vars}(e) \subseteq \Gamma \\ \Gamma; \{ x \mapsto e * P \} \rightsquigarrow \{ x \mapsto e * Q \} \mid c \end{array}}{\Gamma; \{ x \mapsto - * P \} \rightsquigarrow \{ x \mapsto e * Q \} \mid *x = e; c} \quad (\text{Write})$$

$\{ x \mapsto a * y \mapsto b \}$

`void swap(loc x, loc y)`

$\{ x \mapsto b * y \mapsto a \}$

$$\{x, y\}; \{x \mapsto a * y \mapsto b\} \rightsquigarrow \{x \mapsto b * y \mapsto a\} \quad | \quad ??$$

$\{x, y, a2\}; \{x \mapsto a2 * y \mapsto b\} \rightsquigarrow \{x \mapsto b * y \mapsto a2\} \mid ??$

$\{x, y\}; \{x \mapsto a * y \mapsto b\} \rightsquigarrow \{x \mapsto b * y \mapsto a\} \mid \text{let } a2 = *x; ??$

(Read)

$\{x, y, a2, b2\}; \{x \mapsto a2 * y \mapsto b2\} \rightsquigarrow \{x \mapsto b2 * y \mapsto a2\} \mid ??$

(Read)

$\{x, y, a2\}; \{x \mapsto a2 * y \mapsto b\} \rightsquigarrow \{x \mapsto b * y \mapsto a2\} \mid \text{let } b2 = *y; ??$

(Read)

$\{x, y\}; \{x \mapsto a * y \mapsto b\} \rightsquigarrow \{x \mapsto b * y \mapsto a\} \mid \text{let } a2 = *x; ??$

$$\{x, y, a2, b2\}; \{x \mapsto b2 * y \mapsto b2\} \rightsquigarrow \{x \mapsto b2 * y \mapsto a2\} \mid ??$$

(Write)

$$\{x, y, a2, b2\}; \{x \mapsto a2 * y \mapsto b2\} \rightsquigarrow \{x \mapsto b2 * y \mapsto a2\} \mid *x = b2; ??$$

(Read)

$$\{x, y, a2\}; \{x \mapsto a2 * y \mapsto b\} \rightsquigarrow \{x \mapsto b * y \mapsto a2\} \mid \text{let } b2 = *y; ??$$

(Read)

$$\{x, y\}; \{x \mapsto a * y \mapsto b\} \rightsquigarrow \{x \mapsto b * y \mapsto a\} \mid \text{let } a2 = *x; ??$$

$\{x, y, a2, b2\}; \{y \mapsto b2\} \rightsquigarrow \{y \mapsto a2\} \mid ??$

(Frame)

$\{x, y, a2, b2\}; \{x \mapsto b2 * y \mapsto b2\} \rightsquigarrow \{x \mapsto b2 * y \mapsto a2\} \mid ??$

(Write)

$\{x, y, a2, b2\}; \{x \mapsto a2 * y \mapsto b2\} \rightsquigarrow \{x \mapsto b2 * y \mapsto a2\} \mid *x = b2; ??$

(Read)

$\{x, y, a2\}; \{x \mapsto a2 * y \mapsto b\} \rightsquigarrow \{x \mapsto b * y \mapsto a2\} \mid \text{let } b2 = *y; ??$

(Read)

$\{x, y\}; \{x \mapsto a * y \mapsto b\} \rightsquigarrow \{x \mapsto b * y \mapsto a\} \mid \text{let } a2 = *x; ??$

$\{x, y, a2, b2\}; \{y \mapsto a2\} \rightsquigarrow \{y \mapsto a2\} \mid ??$

(Write)

$\{x, y, a2, b2\}; \{y \mapsto b2\} \rightsquigarrow \{y \mapsto a2\} \mid *y = a2; ??$

(Frame)

$\{x, y, a2, b2\}; \{x \mapsto b2 * y \mapsto b2\} \rightsquigarrow \{x \mapsto b2 * y \mapsto a2\} \mid ??$

(Write)

$\{x, y, a2, b2\}; \{x \mapsto a2 * y \mapsto b2\} \rightsquigarrow \{x \mapsto b2 * y \mapsto a2\} \mid *x = b2; ??$

(Read)

$\{x, y, a2\}; \{x \mapsto a2 * y \mapsto b\} \rightsquigarrow \{x \mapsto b * y \mapsto a2\} \mid \text{let } b2 = *y; ??$

(Read)

$\{x, y\}; \{x \mapsto a * y \mapsto b\} \rightsquigarrow \{x \mapsto b * y \mapsto a\} \mid \text{let } a2 = *x; ??$

$\{x, y, a2, b2\}; \{ \text{emp} \} \rightsquigarrow \{ \text{emp} \} \mid ??$

(Frame)

$\{x, y, a2, b2\}; \{ y \mapsto a2 \} \rightsquigarrow \{ y \mapsto a2 \} \mid ??$

(Write)

$\{x, y, a2, b2\}; \{ y \mapsto b2 \} \rightsquigarrow \{ y \mapsto a2 \} \mid *y = a2; ??$

(Frame)

$\{x, y, a2, b2\}; \{ x \mapsto b2 * y \mapsto b2 \} \rightsquigarrow \{ x \mapsto b2 * y \mapsto a2 \} \mid ??$

(Write)

$\{x, y, a2, b2\}; \{ x \mapsto a2 * y \mapsto b2 \} \rightsquigarrow \{ x \mapsto b2 * y \mapsto a2 \} \mid *x = b2; ??$

(Read)

$\{x, y, a2\}; \{ x \mapsto a2 * y \mapsto b \} \rightsquigarrow \{ x \mapsto b * y \mapsto a2 \} \mid \text{let } b2 = *y; ??$

(Read)

$\{x, y\}; \{ x \mapsto a * y \mapsto b \} \rightsquigarrow \{ x \mapsto b * y \mapsto a \} \mid \text{let } a2 = *x; ??$

(Emp)

$\{x, y, a2, b2\}; \{ \text{emp} \} \rightsquigarrow \{ \text{emp} \} \mid \text{skip}$

(Frame)

$\{x, y, a2, b2\}; \{ y \mapsto a2 \} \rightsquigarrow \{ y \mapsto a2 \} \mid ??$

(Write)

$\{x, y, a2, b2\}; \{ y \mapsto b2 \} \rightsquigarrow \{ y \mapsto a2 \} \mid \boxed{*y = a2;} ??$

(Frame)

$\{x, y, a2, b2\}; \{ x \mapsto b2 * y \mapsto b2 \} \rightsquigarrow \{ x \mapsto b2 * y \mapsto a2 \} \mid ??$

(Write)

$\{x, y, a2, b2\}; \{ x \mapsto a2 * y \mapsto b2 \} \rightsquigarrow \{ x \mapsto b2 * y \mapsto a2 \} \mid \boxed{*x = b2;} ??$

(Read)

$\{x, y, a2\}; \{ x \mapsto a2 * y \mapsto b \} \rightsquigarrow \{ x \mapsto b * y \mapsto a2 \} \mid \boxed{\text{let } b2 = *y;} ??$

(Read)

$\{x, y\}; \{ x \mapsto a * y \mapsto b \} \rightsquigarrow \{ x \mapsto b * y \mapsto a \} \mid \boxed{\text{let } a2 = *x;} ??$

```
void swap(loc x, loc y) {  
    let a2 = *x;  
    let b2 = *y;  
    *x = b2;  
    *y = a2;  
}
```

What Can be Synthesised

<i>Data Structure</i>	<i>Id</i>	<i>Description</i>	<i>Proc</i>	<i>Stmt</i>	<i>Code/Spec</i>	<i>Time</i>
Integers	1	swap two	1	4	1.0x	0.2
	2	min of two ¹	1	3	1.1x	0.8
Singly Linked List	3	length ²	1	6	1.4x	0.4
	4	max ²	1	11	1.9x	3.0
	5	min ²	1	11	1.9x	2.9
	6	singleton ¹	1	4	0.9x	0.2
	7	deallocate	1	4	5.5x	0.2
	8	initialize	1	4	1.6x	0.4
	9	copy ³	1	11	2.7x	0.6
	10	append ³	1	6	1.1x	0.4
	11	delete ³	1	12	2.6x	1.2
	12	deallocate two	2	9	6.2x	0.2
	13	append three	2	14	2.3x	1.0
	14	non-destructive append	2	21	3.0x	8.0
	15	union	2	23	5.5x	4.3
	16	intersection ⁴	3	32	7.0x	101.1
	17	difference ⁴	2	21	5.1x	4.7
	18	deduplicate ⁴	2	22	7.3x	1.8
Sorted list	19	prepend ²	1	4	0.4x	0.2
	20	insert ²	1	19	3.1x	1.0
	21	insertion sort ²	1	7	1.2x	0.7
	22	sort ⁴	2	13	4.9x	1.0
	23	reverse ⁴	2	11	4.0x	0.7
	24	merge ²	2	30	4.4x	55.6
Doubly Linked List	25	singleton ¹	1	5	1.1x	0.2
	26	copy	1	22	4.3x	7.2
	27	append ³	1	10	1.6x	1.7
	28	delete ³	1	19	3.7x	3.4
	29	single to double	1	23	6.0x	0.7

What Can be Synthesised

<i>Data Structure</i>	<i>Id</i>	<i>Description</i>	<i>Proc Stmt Code/Spec</i>			<i>Time</i>
Doubly Linked List	25	singleton ¹	1	5	1.1x	0.2
	26	copy	1	22	4.3x	7.2
	27	append ³	1	10	1.6x	1.7
	28	delete ³	1	19	3.7x	3.4
	29	single to double	1	23	6.0x	0.7
List of Lists	30	deallocate	2	11	10.7x	0.2
	31	flatten ⁴	2	17	4.4x	0.6
	32	length ⁵	2	21	5.5x	22.8
Binary Tree	33	size	1	9	2.5x	0.4
	34	deallocate	1	6	8.0x	0.2
	35	deallocate two	1	16	11.8x	0.4
	36	copy	1	16	3.8x	2.5
	37	flatten w/append	1	17	4.8x	0.4
	38	flatten w/acc	1	12	2.1x	0.6
	39	flatten	2	23	7.1x	1.5
	40	flatten to dll in place	2	15	9.6x	11.3
	41	flatten to dll w/null ⁵	2	17	11.2x	106.1
BST	42	insert ²	1	19	2.8x	14.6
	43	rotate left ²	1	5	0.2x	6.2
	44	rotate right ²	1	5	0.2x	4.9
	45	find min ⁵	1	11	1.4x	66.3
	46	find max ⁵	1	18	2.2x	58.0
	47	delete root ²	1	18	1.3x	13.9
	48	from list ⁴	2	27	5.7x	10.0
	49	to sorted list ⁴	3	32	7.7x	20.8
Rose Tree	50	deallocate	2	9	12.0x	0.2
	51	flatten	3	25	8.0x	11.0
	52	copy ⁵	2	32	7.9x	-
Packed Tree	53	pack ⁵	1	16	1.6x	-
	54	unpack ⁵	1	23	2.9x	21.0

More Program Synthesis

Excel® FlashFill

Excel sees patterns and shows a preview

	A	B	C	D	E	F	G
1	Name	First	Last				
2	Ned Lanning	Ned					
3	Margo Hendrix	Margo					
4	Dianne Pugh	Dianne					
5	Earlene McCarty	Earlene					
6	Jon Voigt	Jon					
7	Mia Arnold	Mia					
8	Jorge Fellows	Jorge					
9	Rose Winters	Rose					
10	Carmela Hahn	Carmela					
11	Denis Horning	Denis					
12	Johnathan Swope	Johnathan					
13	Delia Cochran	Delia					
14	Marguerite Cervantes	Marguerite					
15	Liliana English	Liliana					
16	Wendy Stephenson	Wendy					

Wrapping Up

Why CS4212?

- We have learned (hopefully):
 - How high-level languages are implemented in machine language
 - (A subset of) Intel x86 architecture
 - (A subset of) LLVM
 - Lexing and parsing
 - Lambda-calculus and its extensions
 - A little about programming language semantics and type systems
 - How to implement a type checker
 - How to implement a program analyser
 - How to implement a program optimiser
 - Practical applications of theory (logic, proofs, automata, graphs, lattices)
 - How to represent complex data structures in memory
 - How to write large working programs in OCaml
 - How to be a better programmer

Where else is this stuff applicable?

- Understanding hardware/software interface
 - Different devices have different instruction sets, programming models
- General programming
 - In C/C++, better understanding of how the compiler works can help you generate better code.
 - Ability to read assembly/LLVM output from compiler
 - Experience with functional programming gives you different ways to think about solving a problem
 - Knowledge of type systems helps you understand type errors in Java, Scala, OCaml, etc.
- Writing domain specific languages
 - lex/yacc very useful for little utilities
 - understanding abstract syntax specification
 - understanding typing rules
 - being able to write *your own* interpreter and compiler

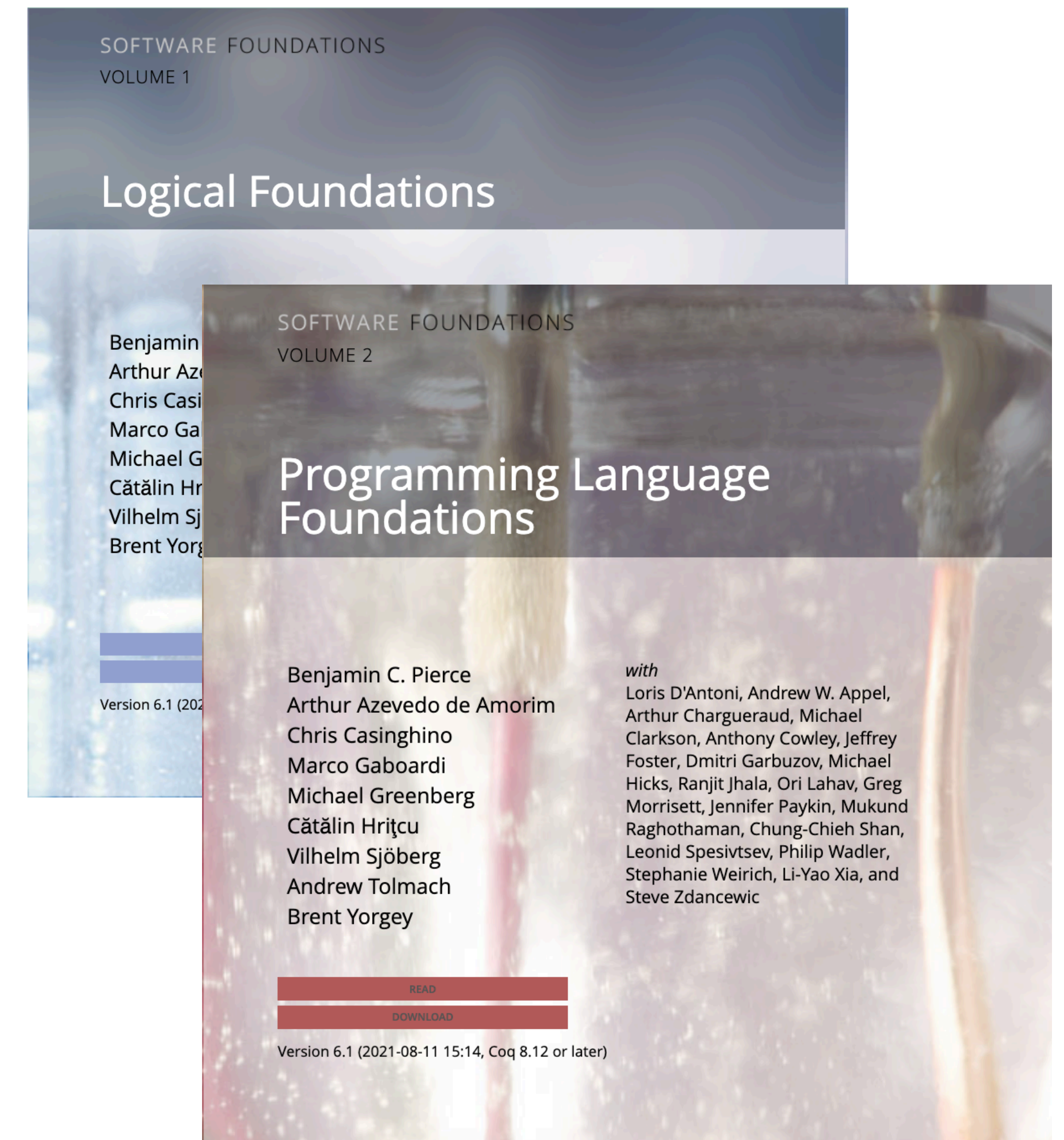
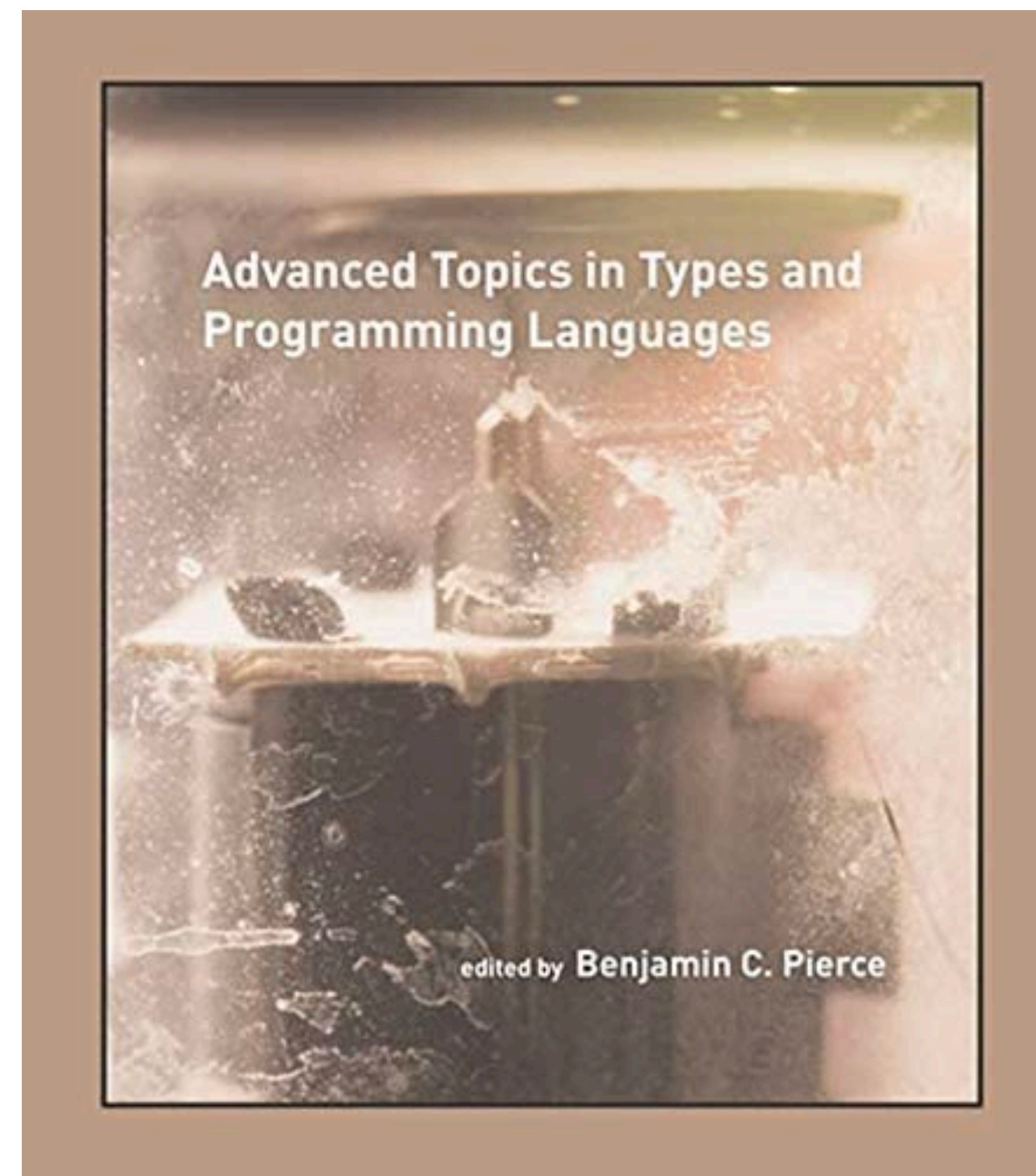
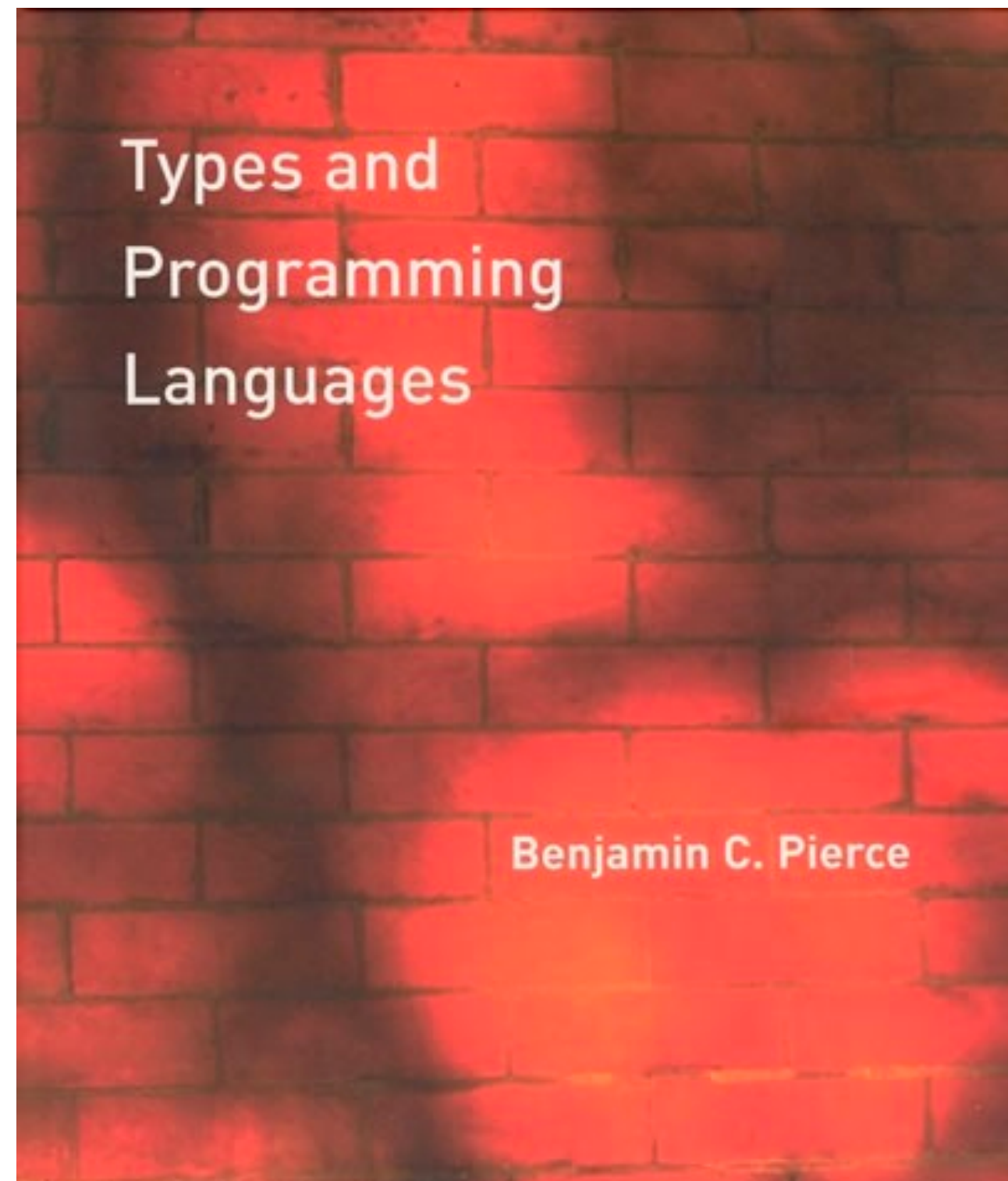
Stuff we didn't cover

- We skipped stuff at every level...
- Concrete syntax/parsing:
 - Much more to the theory of parsing... LR(*)
 - Good syntax is art not science!
- Source language features:
 - Exceptions, advanced type systems, type inference, dependent types, concurrency
- Intermediate languages:
 - Intermediate language design, bytecode, bytecode interpreters, just-in-time compilation (JIT)
- Compilation:
 - Continuation-passing transformation, analyses for SSA, compiling OO classes, lambda-lifting, closure conversion
- Analysis and Optimisations:
 - Abstract interpretation, cache optimization, instruction selection/optimization
- Runtime support:
 - memory management, garbage collection

Where to go from here?

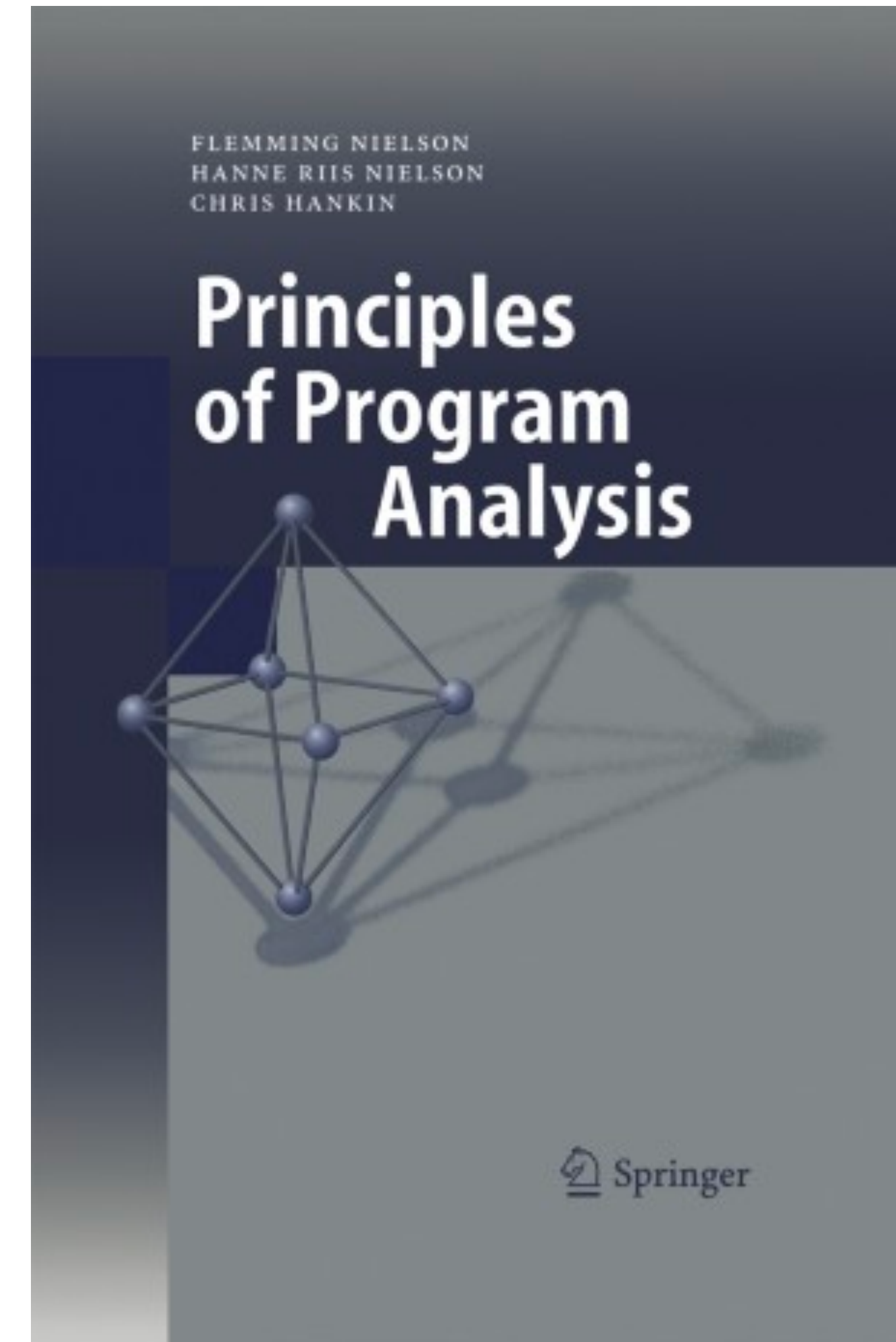
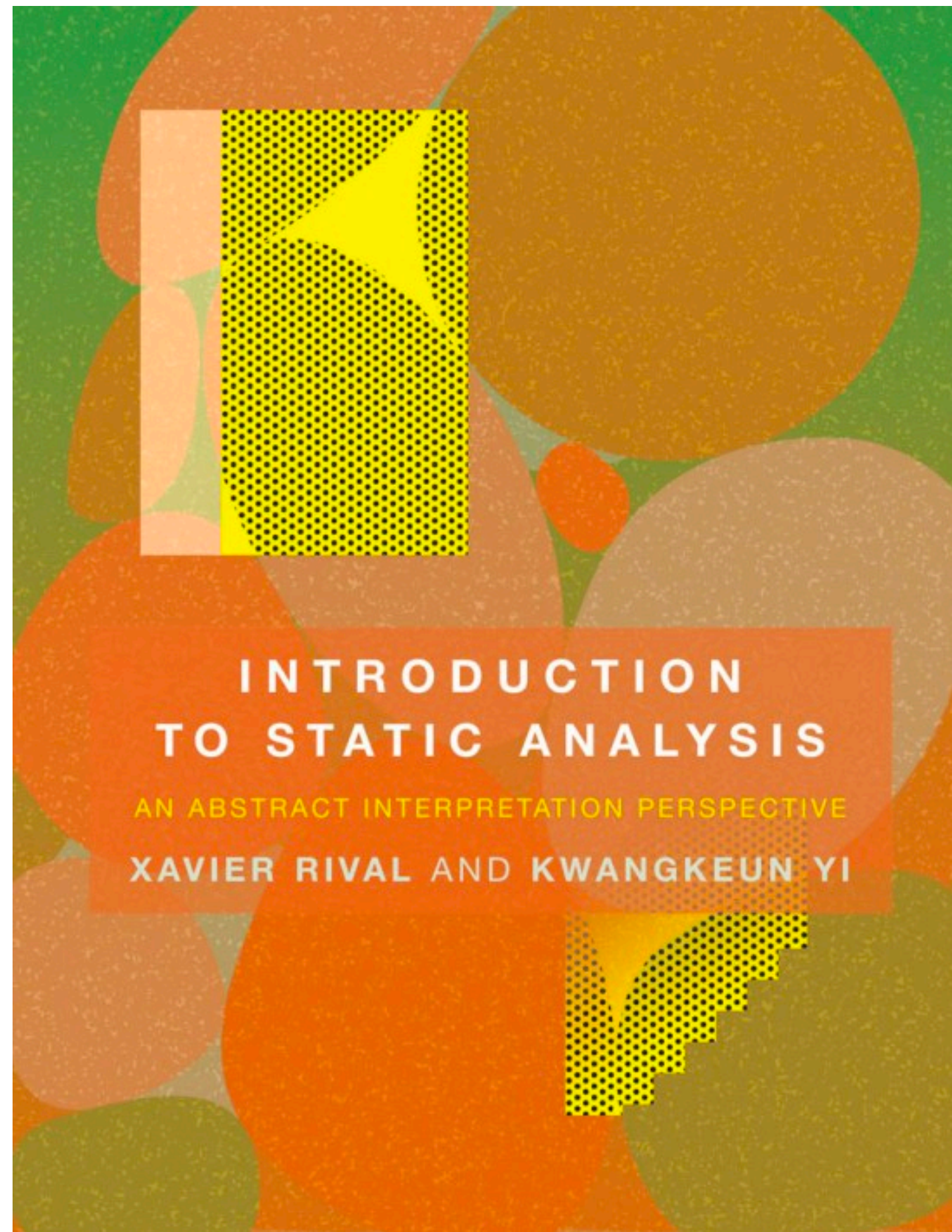
- Conferences (proceedings available on the web):
 - Programming Language Design and Implementation (PLDI)
 - Principles of Programming Languages (POPL)
 - Object Oriented Programming Systems, Languages & Applications (OOPSLA)
 - International Conference on Functional Programming (ICFP)
- Programming Language Mentoring Workshops (PLMW)
 - Affiliated with POPL/PLDI/OOPSLA/ICFP
- Technologies / Open Source Projects
 - Yacc, lex, bison, flex, ...
 - LLVM – low level virtual machine
 - Java virtual machine (JVM), Microsoft's Common Language Runtime (CLR)
 - Languages: OCaml, F#, Haskell, Scala, Go, Rust, ... Coq, Agda, ...?

Further Reading - Types

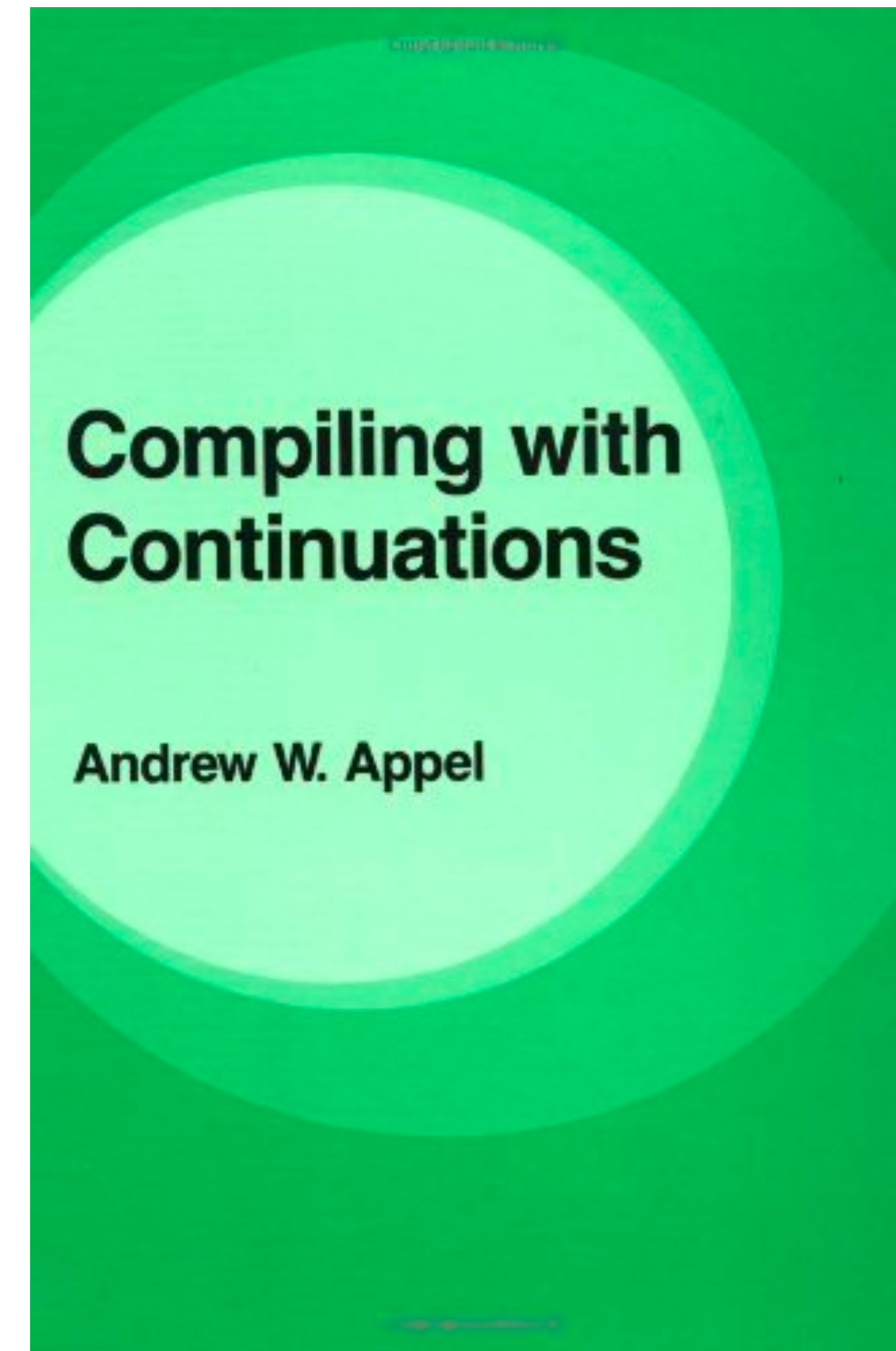


<https://softwarefoundations.cis.upenn.edu/>

Further Reading - Analysis



Further Reading - Compilation Techniques



PL Classes at NUS School of Computing

- CS4215: Programming Language Implementation
 - semantics, type systems, automatic memory management, dynamic linking and just-in-time compilation, as features of modern execution systems
- CS5218: Principles and Practice of Program Analysis
 - foundations of static program analysis, abstract interpretation, lattice theory, analysis of higher-order languages
- CS6217: Topics in Programming Languages & Software Engineering
 - different topics every year, this year: deductive program verification
 - see <https://ilyasergey.net/CS6217/>

PL & Compiler Research @ NUS SoC

nus-plse.github.io

or google “NUS PLSE”



The End

Thanks!