# Intermediate Representations (IRs)

Martin Kellogg

#### Agenda

- quiz
- what is an IR?
  - $\circ$   $\,$  a taxonomy and toolbox  $\,$
- three-address code ("TAC", relevant to PA3c2)
- single static assignment form ("SSA", maybe relevant to other PAs...)

#### Agenda

- quiz
- what is an IR?
  - $\circ$   $\,$  a taxonomy and toolbox  $\,$
- three-address code ("TAC", relevant to PA3c2)
- single static assignment form ("SSA", maybe relevant to other PAs...)

#### Agenda

- quiz
- what is an IR?
  - $\circ$   $\,$  a taxonomy and toolbox  $\,$
- three-address code ("TAC", relevant to PA3c2)
- single static assignment form ("SSA", maybe relevant to other PAs...)

**Definition:** an *intermediate representation* (*IR*) is any internal data structure that a compiler uses for the facts that it derives about the program.

• this definition is intentionally broad, and it includes data structures we've already talked about, including ASTs, CFGs, the implementation/parent/class maps from PA2, etc.

- this definition is intentionally broad, and it includes data structures we've already talked about, including ASTs, CFGs, the implementation/parent/class maps from PA2, etc.
- the IR is generally the canonical form that the compiler reasons about (i.e., the compiler discards the source code)

- this definition is intentionally broad, and it includes data structures we've already talked about, including ASTs, CFGs, the implementation/parent/class maps from PA2, etc.
- the IR is generally the canonical form that the compiler reasons about (i.e., the compiler discards the source code)
  - this makes IR choice important facts that are not in the IR aren't available to the compiler!

- this definition is intentionally broad, and it includes data structures we've already talked about, including ASTs, CFGs, the implementation/parent/class maps from PA2, etc.
- the IR is generally the canonical form that the compiler reasons about (i.e., the compiler discards the source code)
  - this makes IR choice important facts that are not in the IR aren't available to the compiler!
- typically, a compiler uses **different IRs** for different tasks

- As many as are useful!
  - I know that's not a very helpful answer...

- As many as are useful!
  - I know that's not a very helpful answer...
- Typically, a compiler will process IRs in stages
  - i.e., it will transform one IR into another, into another, etc., before eventually emitting code

- As many as are useful!
  - I know that's not a very helpful answer...
- Typically, a compiler will process IRs in stages
  - i.e., it will transform one IR into another, into another, etc., before eventually emitting code
- Different IRs are **better suited to different tasks**

- As many as are useful!
  - I know that's not a very helpful answer...
- Typically, a compiler will process IRs in stages
  - i.e., it will transform one IR into another, into another, etc., before eventually emitting code
- Different IRs are **better suited to different tasks** 
  - for example, a CFG is useful for abstract interpretation or dataflow analysis...

- As many as are useful!
  - I know that's not a very helpful answer...
- Typically, a compiler will process IRs in stages
  - i.e., it will transform one IR into another, into another, etc., before eventually emitting code
- Different IRs are **better suited to different tasks** 
  - for example, a CFG is useful for abstract interpretation or dataflow analysis...
  - ...but would you use it for peephole optimizations? No, you want something much closer to the target assembly

- As many as are useful!
  - An example of a *peephole* I know that's not a very helpful optimization is removing a
- Typically, a compiler will process redundant store/load pair • i.e., it will transform one IR int that access the same address before eventually emitting code
- Different IRs are **better suited to different tasks** 
  - for example, a CFG is useful for abstract interpretation or dataflow analysis...
  - ...but would you use it for peephole optimizations? No, you want something much closer to the target assembly

• The design space of IRs is very large

- The design space of IRs is very large
  - $\circ$   $\,$  In this lecture, I will survey some common ones

- The design space of IRs is very large
  - In this lecture, I will survey some common ones
  - You are welcome to choose to use any or all of these in your compiler (or invent your own)

- The design space of IRs is very large
  - In this lecture, I will survey some common ones
  - You are welcome to choose to use any or all of these in your compiler (or invent your own)
- IR design decisions impact the whole compiler

- The design space of IRs is very large
  - In this lecture, I will survey some common ones
  - You are welcome to choose to use any or all of these in your compiler (or invent your own)
- IR design decisions impact the whole compiler
  - desirable properties:

- The design space of IRs is very large
  - In this lecture, I will survey some common ones
  - You are welcome to choose to use any or all of these in your compiler (or invent your own)
- IR design decisions impact the whole compiler
  - desirable properties:
    - easy to generate
    - easy to manipulate
    - expressive
    - appropriate level of abstraction

- The design space of IRs is very large
  - In this lecture, I will survey some common ones
  - You are welcome to choose to use any or all of these in your compiler (or invent your own)
- IR design decisions impact the whole compiler
  - desirable properties:
    - easy to generate
    - easy to manipulate
    - expressive
    - appropriate level of abstraction
  - tradeoffs between these motivate different IRs at different compilation stages

• A non-exhaustive list of ways that IRs vary:

- A non-exhaustive list of ways that IRs vary:
  - Structure

- A non-exhaustive list of ways that IRs vary:
  - Structure
    - Graphical (trees, graphs, etc.)

- A non-exhaustive list of ways that IRs vary:
  - Structure
    - Graphical (trees, graphs, etc.)
    - Linear (code for some abstract machine)

- A non-exhaustive list of ways that IRs vary:
  - Structure
    - Graphical (trees, graphs, etc.)
    - Linear (code for some abstract machine)
    - Hybrids are common (e.g., control-flow graphs whose nodes are basic blocks of linear code)

- A non-exhaustive list of ways that IRs vary:
  - Structure
    - Graphical (trees, graphs, etc.)
    - Linear (code for some abstract machine)
    - Hybrids are common (e.g., control-flow graphs whose nodes are basic blocks of linear code)
  - Abstraction Level

- A non-exhaustive list of ways that IRs vary:
  - Structure
    - Graphical (trees, graphs, etc.)
    - Linear (code for some abstract machine)
    - Hybrids are common (e.g., control-flow graphs whose nodes are basic blocks of linear code)
  - Abstraction Level
    - High-level, near to source language

- A non-exhaustive list of ways that IRs vary:
  - Structure
    - Graphical (trees, graphs, etc.)
    - Linear (code for some abstract machine)
    - Hybrids are common (e.g., control-flow graphs whose nodes are basic blocks of linear code)
  - Abstraction Level
    - High-level, near to source language
    - Low-level, closer to machine (exposes more details to compiler)

- A non-exhaustive list of ways that IRs vary:
  - Structure
    - Graphical (trees, graphs, etc.)
    - Linear (code for some abstract machine)
    - Hybrids are common (e.g., control-flow graphs whose nodes are basic blocks of linear code)
  - Abstraction Level
    - High-level, near to source language
    - Low-level, closer to machine (exposes more details to compiler)
  - Naming conventions

- A non-exhaustive list of w
  - **Structure** 
    - Graphical (trees,
    - Linear (code for s
- Think of our discussion today as a toolbox of various well-known IRs. When building your compiler, pick and choose which ones are useful for Hybrids are comn each task

nodes are basic blocks of linear code

- Abstraction Level
  - High-level, near to source language
  - Low-level, closer to machine (exposes more details to compiler)
- Naming conventions Ο

#### Example: Representing an Array Reference

#### Example: Representing an Array Reference

source code: A[i,j]

#### Example: Representing an Array Reference

source code: A[i,j]

graphical IR (tree):


#### Example: Representing an Array Reference

source code: A[i,j]

graphical IR (tree):



high-level linear IR: t<sub>1</sub> <- A[i,j]

#### **Example: Representing an Array Reference**

source code: A[i,j]

graphical IR (tree):



high-level linear IR: t<sub>1</sub> <- A[i,j]

low-level	loadI 1 => r1
linear IR:	$sub rj, r1 \Rightarrow r2$
	loadI 10 => r3
	mult $r2, r3 => r4$
	sub <i>ri</i> ,r1 => r5
	add r4,r5 => r6
	loadI @A => r7
	add r7,r6 => r8
÷ 1	load r8 => r9

• Key design decision: how much **detail** to expose

- Key design decision: how much **detail** to expose
  - Affects possibility and profitability of various optimizations

- Key design decision: how much **detail** to expose
  - Affects possibility and profitability of various optimizations
  - Depends on compiler phase:

- Key design decision: how much **detail** to expose
  - Affects possibility and profitability of various optimizations
  - Depends on compiler phase:
    - Some semantic analysis & optimizations are easier with high-level IRs close to the source code (e.g., an AST)

- Key design decision: how much **detail** to expose
  - Affects possibility and profitability of various optimizations
  - Depends on compiler phase:
    - Some semantic analysis & optimizations are easier with high-level IRs close to the source code (e.g., an AST)
    - Low-level usually preferred for other optimizations, register allocation, code generation, etc.

- Key design decision: how much **detail** to expose
  - Affects possibility and profitability of various optimizations
  - Depends on compiler phase:
    - Some semantic analysis & optimizations are easier with high-level IRs close to the source code (e.g., an AST)
    - Low-level usually preferred for other optimizations, register allocation, code generation, etc.
  - Structural (graphical) IRs are typically fairly high-level but are also used for low-level

- Key design decision: how much **detail** to expose
  - Affects possibility and profitability of various optimizations
  - Depends on compiler phase:
    - Some semantic analysis & optimizations are easier with high-level IRs close to the source code (e.g., an AST)
    - Low-level usually preferred for other optimizations, register allocation, code generation, etc.
  - Structural (graphical) IRs are typically fairly high-level but are also used for low-level
  - Linear IRs are typically low-level

- Key design decision: how much **detail** to expose
  - Affects possibility and profitability of various optimizations
  - Depends on compiler phase:
    - Some semantic analysis & optimizations are easier with high-level IRs close to the source code (e.g., an AST)
    - Low-level usually preferred for other optimizations, register allocation, code generation, etc.
  - Structural (graphical) IRs are typically fairly high-level but are also used for low-level
  - Linear IRs are typically low-level
  - But these generalizations don't always hold

- Key design decision: how
  - Affects possibility ar O
  - Depends on compile
    - Some semantic a high-level IRs clo

One view: **source code** is *just another IR* that we happen to expose to programmers

- high-level of abstraction, linear
- Low-level usually register allocation, code generation, etc.
- Structural (graphical) IRs are typically fairly high-level but are also used for low-level
- Linear IRs are typically low-level
- But these generalizations don't always hold

• A graphical IR is any IR that is represented as a graph (or tree, flowchart, or other "graphical" structure)

- A graphical IR is any IR that is represented as a graph (or tree, flowchart, or other "graphical" structure)
- Nodes and edges typically reflect structure of the program
  - E.g., control flow, data dependence, caller/callee

- A graphical IR is any IR that is represented as a graph (or tree, flowchart, or other "graphical" structure)
- Nodes and edges typically reflect structure of the program
  E.g., control flow, data dependence, caller/callee
- May be large (especially syntax trees)

- A graphical IR is any IR that is represented as a graph (or tree, flowchart, or other "graphical" structure)
- Nodes and edges typically reflect structure of the program
  E.g., control flow, data dependence, caller/callee
- May be large (especially syntax trees)
- High-level examples: syntax trees, directed-acyclic graphs (DAGs)
  - Common in early phases of compilers

- A graphical IR is any IR that is represented as a graph (or tree, flowchart, or other "graphical" structure)
- Nodes and edges typically reflect structure of the program
  E.g., control flow, data dependence, caller/callee
- May be large (especially syntax trees)
- High-level examples: syntax trees, *directed-acyclic graphs* (DAGs)
  Common in early phases of compilers
- Other examples: control flow graphs, *data dependency graphs*, and *call graphs* 
  - Often used in semantic analysis, optimization, and code generation

• Can use a directed acyclic graph to store a variation on the AST

- Can use a directed acyclic graph to store a variation on the AST
  - to capture shared substructures

- Can use a directed acyclic graph to store a variation on the AST
  - $\circ$  to capture shared substructures
- Example: (a \* 2) + ((a \* 2) \* b)



- Can use a directed acyclic graph to store a variation on the AST
  - to capture shared substructures
- Example: (a \* 2) + ((a \* 2) \* b)
- Pros: saves space, exposes redundant sub-expressions



- Can use a directed acyclic graph to store a variation on the AST
  - to capture shared substructures
- Example: (a \* 2) + ((a \* 2) \* b)
- Pros: saves space, exposes redundant sub-expressions
  - why might it be useful to expose redundant sub-expressions?



- Can use a directed acyclic graph to store a variation on the AST
  - to capture shared substructures
- Example: (a \* 2) + ((a \* 2) \* b)
- Pros: saves space, exposes redundant sub-expressions
  - why might it be useful to expose redundant sub-expressions?
- Cons: less flexibility if part of tree should be changed



- Can use a directed acyclic graph to store a variation on the AST
  - to capture shared substructures
- Example: (a \* 2) + ((a \* 2) \* b)
- Pros: saves space, exposes redundant sub-expressions
  - why might it be useful to expose redundant sub-expressions?
- Cons: less flexibility if part of tree should be changed
- Ask me about egraphs in OH



• In a *data dependency graph*, an edge between a pair of nodes indicates that they reference common data

- In a *data dependency graph*, an edge between a pair of nodes indicates that they reference common data
- Examples:

- In a *data dependency graph*, an edge between a pair of nodes indicates that they reference common data
- Examples:
  - Block A defines x then B reads it (**RAW** *read after write*)

- In a *data dependency graph*, an edge between a pair of nodes indicates that they reference common data
- Examples:
  - Block A defines x then B reads it (RAW read after write)
  - Block A reads x then B writes it (WAR "anti- dependence")

- In a *data dependency graph*, an edge between a pair of nodes indicates that they reference common data
- Examples:
  - Block A defines x then B reads it (RAW read after write)
  - Block A reads x then B writes it (WAR "anti- dependence")
  - Blocks A and B both write x (WAW)
    - order of blocks must reflect original program semantics

- In a *data dependency graph*, an edge between a pair of nodes indicates that they reference common data
- Examples:
  - Block A defines x then B reads it (RAW read after write)
  - Block A reads x then B writes it (WAR "anti- dependence")
  - Blocks A and B both write x (WAW)
    - order of blocks must reflect original program semantics
- These dependencies restrict what reorderings the compiler can do

- In a *data dependency graph*, an edge between a pair of nodes indicates that they reference common data
- Examples:
  - Block A defines x then B reads it (RAW read after write)
  - Block A reads x then B writes it (WAR "anti- dependence")
  - Blocks A and B both write x (WAW)
    - order of blocks must reflect original program semantics
- These dependencies restrict what reorderings the compiler can do
- Data dependency graph is most often used in conjunction with another IR to facilitate optimizations, but it has other uses too...

#### Call Graphs

## **Call Graphs**

• A *call graph* represents the runtime transfers of control between procedures

# Call Graphs

- A *call graph* represents the runtime transfers of control between procedures
  - one node for each procedure and one edge for each distinct procedure call site
- A *call graph* represents the runtime transfers of control between procedures
  - one node for each procedure and one edge for each distinct procedure call site
  - e.g., if the code calls *q* from three textually distinct sites in *p*; the call graph has three edges (*p*, *q*), one for each call site

- A *call graph* represents the runtime transfers of control between procedures
  - one node for each procedure and one edge for each distinct procedure call site
  - e.g., if the code calls q from three textually distinct sites in p;
     the call graph has three edges (p, q), one for each call site
- Call graphs are useful if you want to do *inter-procedural* analysis

- A *call graph* represents the runtime transfers of control between procedures
  - one node for each procedure and one edge for each distinct procedure call site
  - e.g., if the code calls *q* from three textually distinct sites in *p*; the call graph has three edges (*p*, *q*), one for each call site
- Call graphs are useful if you want to do *inter-procedural* analysis
   that is, analysis that requires you to reason about more than one procedure at the same time

- A *call graph* represents the runtime transfers of control between procedures
  - one node for each procedure and one edge for each distinct procedure call site
  - e.g., if the code calls *q* from three textually distinct sites in *p*; the call graph has three edges (*p*, *q*), one for each call site
- Call graphs are useful if you want to do *inter-procedural* analysis
  - that is, analysis that requires you to reason about more than one procedure at the same time
  - during codegen, you *probably* don't need to do this, but it might be useful for optimization

#### **Trivia Break: International Relations**

This ancient Athenian historian and general chronicled the fifth-century BC war between Athens and Sparta in his History of the Peloponnesian War. He is considered by some to be the first modern political theorist, because of his claims to have applied strict standards of impartiality and evidence-gathering and analysis of cause and effect, without reference to intervention by the gods. The Melian dialogue from his work is regarded as a seminal text of international relations theory.

• Pseudo-code for some abstract machine

- Pseudo-code for some abstract machine
  - "Linear" because, like source code, it has a textual structure

- Pseudo-code for some abstract machine
  - "Linear" because, like source code, it has a textual structure
  - $\circ$  Level of abstraction varies

- Pseudo-code for some abstract machine
  - "Linear" because, like source code, it has a textual structure
  - Level of abstraction varies
- Simple, compact data structures
  - Commonly used: arrays, linked structures

- Pseudo-code for some abstract machine
  - "Linear" because, like source code, it has a textual structure
  - Level of abstraction varies
- Simple, compact data structures
  - Commonly used: arrays, linked structures
- Examples:

- Pseudo-code for some abstract machine
  - "Linear" because, like source code, it has a textual structure
  - Level of abstraction varies
- Simple, compact data structures
  - Commonly used: arrays, linked structures
- Examples:
  - three-address code ("TAC"), which we'll see some examples of in a few minutes (and which you must generate for PA3c2

- Pseudo-code for some abstract machine
  - "Linear" because, like source code, it has a textual structure
  - Level of abstraction varies
- Simple, compact data structures
  - Commonly used: arrays, linked structures
- Examples:
  - three-address code ("TAC"), which we'll see some examples of in a few minutes (and which you must generate for PA3c2
  - stack machine code, which we'll see a long example of in a later lecture

- Pseudo-code for some abstract machine
  - "Linear" because, like source code, it has a textual structure
  - Level of abstraction varies
- Simple, compact data structures
  - Commonly used: arrays, linked structures
- Examples:
  - three-address code ("TAC"), which we'll see some examples of in a few minutes (and which you must generate for PA3c2
  - stack machine code, which we'll see a long example of in a later lecture
  - o single static assignment form, which we'll see briefly today too

• Linear IRs can be close to the source language, very low-level, or somewhere in between.

- Linear IRs can be close to the source language, very low-level, or somewhere in between.
- Consider, for example, linear IRs for C array reference

#### a[i][j+2]

- Linear IRs can be close to the source language, very low-level, or somewhere in between.
- Consider, for example, linear IRs for C array reference

#### a[i][j+2]

• A high-level linear IR might represent this very similarly to source code:

#### t1 <- a[i,j+2]

• A "medium level" IR:

```
t1 <- j + 2
t2 <- i * 20
t3 <- t1 + t2
t4 <- 4 * t3
t5 <- addr a
t6 <- t5 + t4
t7 <- *t6</pre>
```

• A "medium level" IR:

```
t1 <- j + 2
t2 <- i * 20
t3 <- t1 + t2
t4 <- 4 * t3
t5 <- addr a
t6 <- t5 + t4
t7 <- *t6</pre>
```

still retains basic symbolic info about variables

- A "medium level" IR:
- t1 <- j + 2
  t2 <- i \* 20
  t3 <- t1 + t2
  t4 <- 4 \* t3
  t5 <- addr a
  t6 <- t5 + t4
  t7 <- \*t6</pre>

still retains basic symbolic info about variables

- A "low level" IR:
- r1 <- [fp-4] r2 <- r1 + 2 r3 <- [fp-8] r4 <- r3 \* 20 r5 <- r4 + r2 r6 <- 4 \* r5 r7 <- fp - 216 f1 <- [r7+r6]</pre>

- A "medium level" IR:
- t1 <- j + 2
  t2 <- i \* 20
  t3 <- t1 + t2
  t4 <- 4 \* t3
  t5 <- addr a
  t6 <- t5 + t4
  t7 <- \*t6</pre>

still retains basic symbolic info about variables

- A "low level" IR:
- r1 <- [fp-4] r2 <- r1 + 2 r3 <- [fp-8] r4 <- r3 \* 20 r5 <- r4 + r2 r6 <- 4 \* r5 r7 <- fp - 216 f1 <- [r7+r6]</pre>

exposes all details of the low-level layout; explicit memory references and calculations

• **High-level**: good for some high-level optimizations, semantic checking; but can't optimize things that are hidden – like address arithmetic for array subscripting

- **High-level**: good for some high-level optimizations, semantic checking; but can't optimize things that are hidden like address arithmetic for array subscripting
- Low-level: need for good code generation and resource utilization in back end but loses some semantic knowledge (e.g., variables, data aggregates, source relationships are usually missing)

- **High-level**: good for some high-level optimizations, semantic checking; but can't optimize things that are hidden like address arithmetic for array subscripting
- Low-level: need for good code generation and resource utilization in back end but loses some semantic knowledge (e.g., variables, data aggregates, source relationships are usually missing)
- Medium-level: more detail but keeps more higher-level semantic information great for machine-independent optimizations. Many (all?) optimizing compilers work at this level

- **High-level**: good for some high-level optimizations, semantic checking; but can't optimize things that are hidden like address arithmetic for array subscripting
- Low-level: need for good code generation and resource utilization in back end but loses some semantic knowledge (e.g., variables, data aggregates, source relationships are usually missing)
- Medium-level: more detail but keeps more higher-level semantic information – great for machine-independent optimizations. Many (all?) optimizing compilers work at this level
- Many compilers use all 3 in different phases

- Originally used for stack-based computers
  - famous example: B5000, ~1961

- Originally used for stack-based computers
  - famous example: B5000, ~1961
- Often used for virtual machines
  - Classic examples: Pascal's pcode and Java bytecode

- Originally used for stack-based computers
  - famous example: B5000, ~1961
- Often used for virtual machines
  - Classic examples: Pascal's pcode and Java bytecode
- Advantages:

- Originally used for stack-based computers
  - famous example: B5000, ~1961
- Often used for virtual machines
  - Classic examples: Pascal's pcode and Java bytecode
- Advantages:
  - **Compact**; mostly 0-address opcodes (great if sent via network)

- Originally used for stack-based computers
  - famous example: B5000, ~1961
- Often used for virtual machines
  - Classic examples: Pascal's pcode and Java bytecode
- Advantages:
  - **Compact**; mostly 0-address opcodes (great if sent via network)
  - Easy to generate; easy to write a front-end compiler, leaving the "heavy lifting" and optimizations to the JIT

- Originally used for stack-based computers
  - famous example: B5000, ~1961
- Often used for virtual machines
  - Classic examples: Pascal's pcode and Java bytecode
- Advantages:
  - Compact; mostly 0-address opcodes (great if sent via network)
  - Easy to generate; easy to write a front-end compiler, leaving the "heavy lifting" and optimizations to the JIT
  - Simple to interpret or compile to machine code

- Originally used for stack-based computers
  - famous example: B5000, ~1961
- Often used for virtual machines
  - Classic examples: Pascal's pcode and Java bytecode
- Advantages:
  - Compact; mostly 0-address opcodes (great if sent via network)
  - Easy to generate; easy to write a front-end compiler, leaving the "heavy lifting" and optimizations to the JIT
  - Simple to interpret or compile to machine code
- Disadvantages:

- Originally used for stack-based computers
  - famous example: B5000, ~1961
- Often used for virtual machines
  - Classic examples: Pascal's pcode and Java bytecode
- Advantages:
  - Compact; mostly 0-address opcodes (great if sent via network)
  - Easy to generate; easy to write a front-end compiler, leaving the "heavy lifting" and optimizations to the JIT
  - Simple to interpret or compile to machine code
- Disadvantages:
  - Somewhat inconvenient/difficult to optimize directly

- Originally used for stack-based computers
  - famous example: B5000, ~1961
- Often used for virtual machines
  - Classic examples: Pascal's pcode and Java bytecode
- Advantages:
  - Compact; mostly 0-address opcodes (great if sent via network)
  - Easy to generate; easy to write a front-end compiler, leaving the "heavy lifting" and optimizations to the JIT
  - Simple to interpret or compile to machine code
- Disadvantages:
  - Somewhat inconvenient/difficult to optimize directly
  - Does not match up with modern chip architectures
| pushaddr  | х |
|-----------|---|
| pushconst | 2 |
| pushval   | n |
| pushval   | m |
| add       |   |
| mult      |   |
| store     |   |

push all operands onto stack	pushaddr x pushconst 2 pushval n pushval m
	add mult store

m	
n	
2	
@x	
?	









- Note compactness:
  - common opcodes just 1 byte wide
  - instructions have 0 or 1 operand

Hypothetical code for x = 2 \* (m + n):



instructions have 0 or 1 operand Ο

Ο

• Combinations of linear and graphical IRs are common

- Combinations of linear and graphical IRs are common
  - for example, a CFG's basic blocks usually contain code in some linear IR (e.g., TAC)

- Combinations of linear and graphical IRs are common
  - for example, a CFG's basic blocks usually contain code in some linear IR (e.g., TAC)
  - we call these *hybrid IRs*

- Combinations of linear and graphical IRs are common
  - for example, a CFG's basic blocks usually contain code in some linear IR (e.g., TAC)
  - we call these *hybrid IRs*
- Level of abstraction varies; you can mix and match based on your needs

- Combinations of linear and graphical IRs are common
  - for example, a CFG's basic blocks usually contain code in some linear IR (e.g., TAC)
  - we call these *hybrid IRs*
- Level of abstraction varies; you can mix and match based on your needs
  - when designing your own compiler's internals, it's okay to be creative: pick the representation that makes your life easiest

• Usual form: x <- y op z

- Usual form: x <- y op z
  - One operator

- Usual form: x <- y op z
  - One operator
  - Maximum of 3 names (thus the name)

- Usual form: x <- y op z
  - One operator
  - Maximum of 3 names (thus the name)
- Eg:x = 2 \* (m + n) becomes t1 <- m + n; t2 <- 2 \* t1; x <- t2

- Usual form: x <- y op z
  - One operator
  - Maximum of 3 names (thus the name)

- Usual form: x <- y op z
  - One operator
  - Maximum of 3 names (thus the name)
- Eg: x = 2 \* (m + n) becomes

t1 <- m + n; t2 <- 2 \* t1; x <- t2

• You may prefer:

add t1, m, n; mul t2, 2, t1; movx, t2

• Invent as many new temp names as needed

- Usual form: x <- y op z
  - One operator
  - Maximum of 3 names (thus the name)
- Eg: x = 2 \* (m + n) becomes

t1 <- m + n; t2 <- 2 \* t1; x <- t2

 $\circ$  You may prefer:

add t1, m, n; mul t2, 2, t1; movx, t2

- Invent as many new temp names as needed
  - "expression temps" don't correspond to any user variables; de-anonymize expressions

- Usual form: x <- y op z
  - One operator
  - Maximum of 3 names (thus the name)
- Eg: x = 2 \* (m + n) becomes

t1 <- m + n; t2 <- 2 \* t1; x <- t2

 $\circ$  You may prefer:

add t1, m, n; mul t2, 2, t1; movx, t2

- Invent as many new temp names as needed
  - "expression temps" don't correspond to any user variables; de-anonymize expressions
- Store in a quadruple: <1hs, rhs1, op, rhs2>

• Advantages:

- Advantages:
  - Resembles code for actual machines

- Advantages:
  - Resembles code for actual machines
  - Explicitly names intermediate results

- Advantages:
  - Resembles code for actual machines
  - Explicitly names intermediate results
  - Compact

- Advantages:
  - Resembles code for actual machines
  - Explicitly names intermediate results
  - Compact
  - Often easy to rearrange

- Advantages:
  - Resembles code for actual machines
  - Explicitly names intermediate results
  - Compact
  - Often easy to rearrange
- Why does PA3c2 require TAC?

- Advantages:
  - Resembles code for actual machines
  - Explicitly names intermediate results
  - Compact
  - Often easy to rearrange
- Why does PA3c2 require TAC?
  - I want you to build *an* IR early. You'll need several before you finish PA3.

- Advantages:
  - Resembles code for actual machines
  - Explicitly names intermediate results
  - Compact
  - Often easy to rearrange
- Why does PA3c2 require TAC?
  - I want you to build an IR early. You'll need several before you finish PA3.
  - I think TAC is basically guaranteed to be useful no matter how you design the rest of your compiler
    - and I had to pick something...

- Advantages:
  - Resembles code for actual machines
  - Explicitly names intermediate results
  - Compact
  - Often easy to rearrange
- Why does PA3c2 require TAC?
  - I want you to build an IR early. You'll need several before you finish PA3.
    The PA3c2 page has a long
  - I think TAC is basically how you design the rel\_including pseudocode to generate it.
    - and I had to pick something...

### Single Static Assignment (SSA)

# Single Static Assignment (SSA)

• IR where each variable has **only one definition** in the program text

## Single Static Assignment (SSA)

- IR where each variable has **only one definition** in the program text
  - A single static definition, but that definition can be in a loop, function, or other code that is executed dynamically many times
# Single Static Assignment (SSA)

- IR where each variable has **only one definition** in the program text
  - A single static definition, but that definition can be in a loop, function, or other code that is executed dynamically many times
- Makes many analyses (and related optimizations) more efficient

# Single Static Assignment (SSA)

- IR where each variable has **only one definition** in the program text
  - A single static definition, but that definition can be in a loop, function, or other code that is executed dynamically many times
- Makes many analyses (and related optimizations) more efficient
- Separates values from memory storage locations

# Single Static Assignment (SSA)

- IR where each variable has **only one definition** in the program text
  - A single static definition, but that definition can be in a loop, function, or other code that is executed dynamically many times
- Makes many analyses (and related optimizations) more efficient
- Separates values from memory storage locations
- Complementary to CFG or data dependency graph
  - better for some things, but cannot do everything

#### SSA: Basic Idea

 Basic Idea: for each original variable v, create a new variable v<sub>n</sub> at the nth definition of the original v. Subsequent uses of v use v<sub>n</sub> until the next definition point. E.g.:

#### SSA: Basic Idea

 Basic Idea: for each original variable v, create a new variable v<sub>n</sub> at the nth definition of the original v. Subsequent uses of v use v<sub>n</sub> until the next definition point. E.g.:

#### Original:

#### SSA: Basic Idea

 Basic Idea: for each original variable v, create a new variable v<sub>n</sub> at the nth definition of the original v. Subsequent uses of v use v<sub>n</sub> until the next definition point. E.g.:

Original:	SSA:
a := x + y	$a_1 := x_0 + y_0$
b := a – 1	$b_1 := a_1 - 1$
a := y + b	$a_2 := y_0 + b_2$
b := x * 4	$b_2 := x_0 * 4$
a := a + b	$a_3^2 := a_2^2 + b_2^2$

• This is fine until we reach a merge point:



• This is fine until we reach a merge point:



• This is fine until we reach a merge point:



• Solution: introduce a *φ*-function ("phi function")

• This is fine until we reach a merge point:



- Solution: introduce a *φ*-function ("phi function")
  - semantics:  $a_3$  is assigned to either  $a_1$  or  $a_2$ , depending on which control flow path us used to reach the  $\varphi$ -function

• It doesn't!

- It doesn't!
- φ-functions **don't actually exist** at run time

- It doesn't!
- φ-functions **don't actually exist** at run time
  - $\circ~$  when we're done using the SSA IR, we translate back out of SSA form, removing all  $\phi$ -functions

- It doesn't!
- φ-functions **don't actually exist** at run time
  - $\circ~$  when we're done using the SSA IR, we translate back out of SSA form, removing all  $\phi$ -functions
    - Basically by adding code to copy all SSA x<sub>i</sub> values to the single, non-SSA variable x

- It doesn't!
- φ-functions **don't actually exist** at run time
  - $\circ~$  when we're done using the SSA IR, we translate back out of SSA form, removing all  $\phi$ -functions
    - Basically by adding code to copy all SSA x<sub>i</sub> values to the single, non-SSA variable x
  - For analysis, all we typically need to know is the connection of uses to definitions – no need to "execute" anything

- It doesn't!
- φ-functions **don't actually exist** at run time
  - $\circ~$  when we're done using the SSA IR, we translate back out of SSA form, removing all  $\phi$ -functions
    - Basically by adding code to copy all SSA x<sub>i</sub> values to the single, non-SSA variable x
  - For analysis, all we typically need to know is the connection of uses to definitions no need to "execute" anything
    - So φ-functions are (only) compile-time bookkeeping

• Could simply add φ-functions for every variable at every join point(!)

- Could simply add φ-functions for every variable at every join point(!)
  - called maximal SSA

- Could simply add φ-functions for every variable at every join point(!)
  - called *maximal SSA*
  - this wastes way too much space and time to be useful in practice

- Could simply add φ-functions for every variable at every join point(!)
  - called *maximal SSA*
  - this wastes way too much space and time to be useful in practice
- Instead, use the *path convergence criterion*: insert a φ-function for variable *a* at program point *z* when:

- Could simply add φ-functions for every variable at every join point(!)
  - called *maximal SSA*
  - this wastes way too much space and time to be useful in practice
- Instead, use the *path convergence criterion*: insert a φ-function for variable *a* at program point *z* when:
  - There are blocks x and y, both containing definitions of a, and x != y

- Could simply add φ-functions for every variable at every join point(!)
  - called *maximal SSA*
  - this wastes way too much space and time to be useful in practice
- Instead, use the *path convergence criterion*: insert a φ-function for variable *a* at program point *z* when:
  - There are blocks x and y, both containing definitions of a, and x != y
  - There are non-empty paths from x to z and from y to z

- Could simply add φ-functions for every variable at every join point(!)
  - called *maximal SSA*
  - this wastes way too much space and time to be useful in practice
- Instead, use the *path convergence criterion*: insert a φ-function for variable *a* at program point *z* when:
  - There are blocks x and y, both containing definitions of a, and x != y
  - There are non-empty paths from x to z and from y to z
  - These paths have no common nodes other than z

- Could simply add φ-functions for every variable at every join point(!)
  - called *maximal SSA*
  - this wastes way too much space and time to be useful in practice
- Instead, use the *path convergence criterion*: insert a  $\varphi$ -function for variable *a* at program  $\varphi$  We'll come back to SSA form when
  - There are blocks x and x != y
    There are provide the set of the
  - There are non-empty cactually do this
  - These paths have no common nodes other than z

#### **Course Announcements**

- PA2 due today!
- PA3c1 (codegen testing) is due on Friday
  - all gas, no brakes
- My OH on Wednesday will be later than usual (4-5 instead of 3:30-4:30), because of a CS faculty meeting until 4
  - might even start a little bit later...