

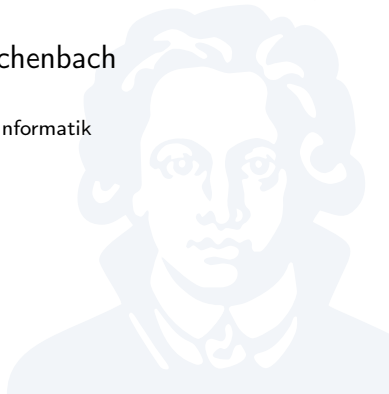
Software Engineering Tools 02

Static Program Analysis

Prof. Dr. Christoph Reichenbach

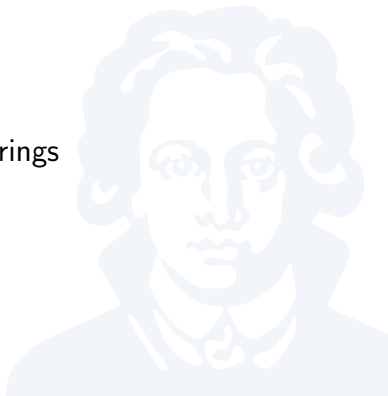
Fachbereich 12 / Institut für Informatik

3. Mai 2014



Why Program Analysis?

- Program understanding
- Finding bugs
- Finding optimisation opportunities
- Checking pre-conditions of refactorings



Program analysis for Optimisation

Example: loop-invariant code motion optimisation

- Move code out of loop if this doesn't alter behaviour

```
int global;
void f(int[] a, int[] b)
{
  for (int i = 0; i < a.length; i++) {
    global = 0; // only need to do this once
    a[i] += b[i];
  }
}
```

- Detect that code is invariant
-
-

Program analysis for Optimisation

Example: loop-invariant code motion optimisation

- Move code out of loop if this doesn't alter behaviour

```
int global;
void f(int[] a, int[] b)
{
  for (int i = 0; i < a.length; i++) {
    a[i] += b[i];
    b[i] = i + Math.atan2(a.length,b.length);
  }
}
```

- Detect that code is invariant
- Invariant code may be *subexpression*
-

Only `Math.atan2(...)` should be moved

Program analysis for Optimisation

Example: loop-invariant code motion optimisation

- Move code out of loop if this doesn't alter behaviour

```
int global;
void f(int[] a, int[] b)
{
    for (int i = 0; i < a.length; i++) {
        a[i] += b[i];
        global = b[7];
    }
}
```

- Detect that code is invariant
- Invariant code may be *subexpression*
- Must consider *pointer aliasing*

Can't move easily if $a = b!$

Today's Topics

- **Type and effect systems**
- Intermediate representations
- **Data-flow analysis**
- **Constraint-based analysis**
- **Abstract Interpretation**
- Analysis Qualities
- Other interesting analyses



Type-and-Effect systems

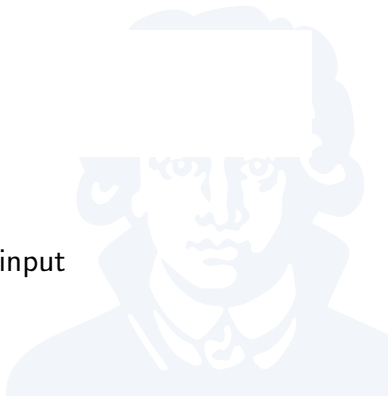
Motivation

- Type analysis can tell us many interesting properties
- How about side effects?

(* Example in Standard ML *)

```
fun f(x) = x*x  
fun g(x) = (print "Foo"; x*x)
```

- Types of `f`, `b`: `int -> int`
- Both compute the square of their input
- `g` *also* produces side effect



Arrows with effect annotations

Effect annotation

$$\begin{array}{c} \downarrow \\ \tau_1 \xrightarrow{\eta} \tau_2 \end{array}$$

- η captures side effects
- Example: $\{\text{io}\} \in F$ for input/output
- We can now distinguish:
 - $f : \text{int} \rightarrow \text{int}$
 - $g : \text{int} \xrightarrow{\{\text{io}\}} \text{int}$

Enrich function arrows with effect annotations

Effect environments

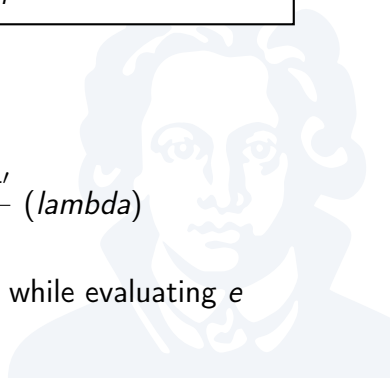
Let F be the set of all effects. We use:

Effect set: $E : 2^F$

- We now type with $\Gamma; E \vdash e$
- Example:

$$\frac{\Gamma + x \mapsto \alpha; E \vdash e : \tau'}{\Gamma; \emptyset \vdash \lambda x. e : \alpha \xrightarrow{E} \tau} \text{ (lambda)}$$

- E captures all effects that happen while evaluating e



Example

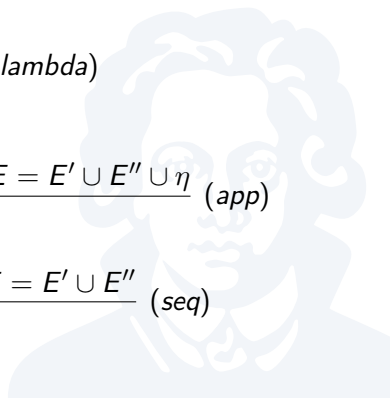
Example for an effect system on a functional language:

$$\frac{\Gamma(n) = \tau}{\Gamma; \emptyset \vdash n : \tau} \text{ (var)}$$

$$\frac{\Gamma + x \mapsto \alpha; E \vdash e : \tau'}{\Gamma; \emptyset \vdash \lambda x. e : \alpha \xrightarrow{E} \tau} \text{ (lambda)}$$

$$\frac{\Gamma; E' \vdash f : \tau_1 \xrightarrow{\eta} \tau_2 \quad \Gamma; E'' \vdash e : \tau_1 \quad E = E' \cup E'' \cup \eta}{\Gamma; E \vdash f e : \tau_2} \text{ (app)}$$

$$\frac{\Gamma; E' \vdash e_1 : \alpha \quad \Gamma; E'' \vdash e_2 : \tau \quad E = E' \cup E''}{\Gamma; E \vdash (e_1; e_2) : \tau} \text{ (seq)}$$

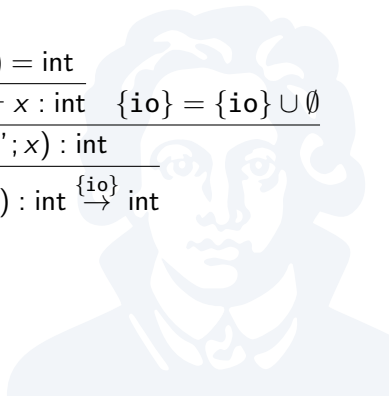


Analysis

$$\frac{\Gamma(\text{print}) = \text{string} \xrightarrow{\{\text{io}\}} \text{unit}}{\Gamma'; \emptyset \vdash \text{print} : \text{string} \xrightarrow{\{\text{io}\}} \text{unit} \quad \dots \quad \frac{\Gamma'(x) = \text{int}}{\Gamma'; \emptyset \vdash x : \text{int}} \quad \frac{\Gamma'; \emptyset \vdash x : \text{int} \quad \{\text{io}\} = \{\text{io}\} \cup \emptyset}{\Gamma'; \{\text{io}\} \vdash (\text{print } \dots; x) : \text{int}}}$$

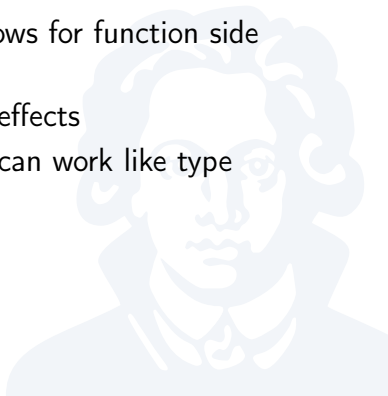
$$\Gamma; \emptyset \vdash (\lambda x. (\text{print } \dots; x)) : \text{int} \xrightarrow{\{\text{io}\}} \text{int}$$

where $\Gamma' = \Gamma + x \mapsto \text{int}$



Summary: Type-and-Effect Systems

- Extension of type systems
- *effect annotations* on function arrows for function side effects
- Effect environments for local side effects
- Not covered here: *effect variables* can work like type variables



Intermediate Representation

Intermediate representations (IRs) simplify compiler construction:

- Encode knowledge needed for many analyses
- Allow compiler backend, optimisers, frontend to evolve separately
- Allow simplifying the language that optimisers/backend operate on
⇒ faster evolution, fewer bugs

Disadvantages:

- Require AST → IR translation
- For *non-optimising* compilers: slow-down, additional complexity

IRs used in practically all “real-world” compilers

Basic Blocks

```
int x = 0;
int y = a;
if (y > x) {
  x = 1;
} else {
  x = 2;
  y = 0;
}
y = x + y;
return y;
```



```
int x = 0;
int y = a;
if (y > x)
```

```
x = 1;
```

```
x = 2;
y = 0;
```

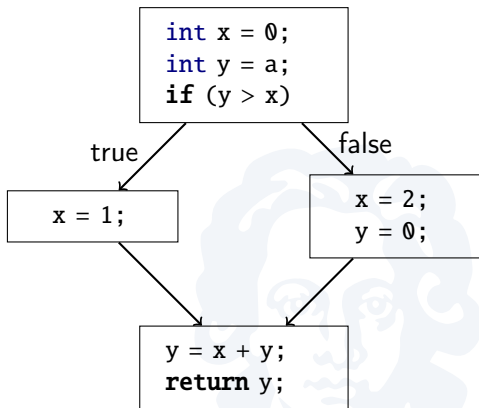
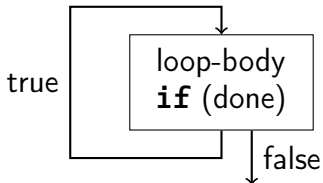
```
y = x + y;
return y;
```

Group operations without intervening control flow

Control-Flow Graphs

Maintain *control edges*

- Labelled: 'true'/false for conditional jumps
- Unlabelled: unconditional jump
- Loops: Back-edges

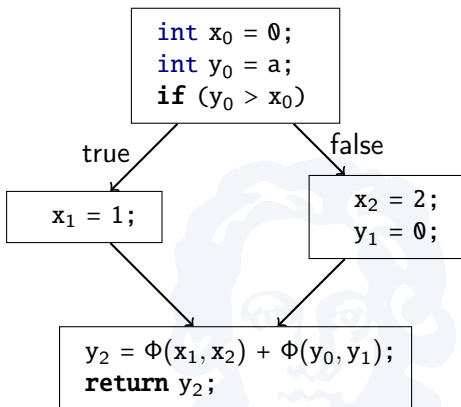


Needed in order for basic blocks to make sense

Static Single-Assignment Form (SSA)

Assign each variable exactly once

- Rename variables to ensure this property
- When different variables flow together, mark with $\Phi(v_1, \dots, v_k)$
- Φ makes dependencies explicit, has trivial semantics (exactly one v_i will be set)

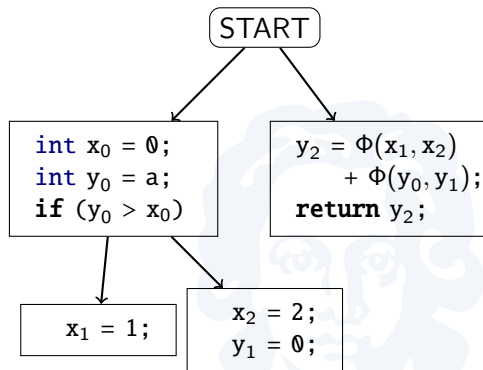


Simplifies or enables many analyses

Control-Dependence Graph

Capture nontrivial dependencies between basic blocks

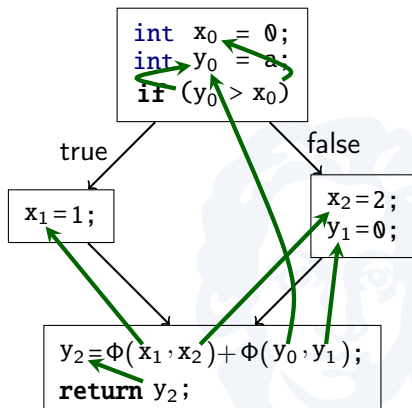
- Arrow $a \rightarrow b$: Block a decides whether b is invoked



Data-Dependence Graph

Capture **dependencies** between reads and assignments

- $a \rightarrow b$: b stores a value that a reads

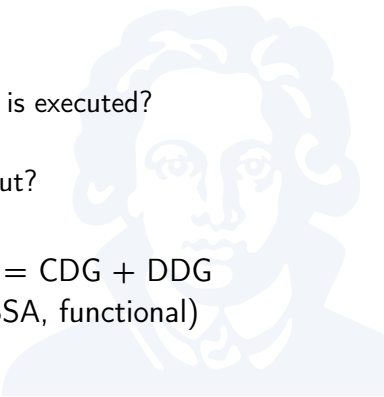


Summary: Intermediate Representations

- Basic Blocks
 - Group operations that always happen together
- SSA (Static Single-Assignment Form)
 - Bind each name exactly once
 - Φ nodes merge multiple inputs
- CDG (Control-Dependence Graph)
 - Who decides whether something is executed?
- DDG (Data-Dependence Graph)
 - Who assigns variables I care about?

Other IRs:

- Program dependence Graph, PDG = CDG + DDG
- A-Normal Form, ANF (similar to SSA, functional)
- Continuation-Passing Form, CPS
(ANF plus jumps, exceptions, ...)

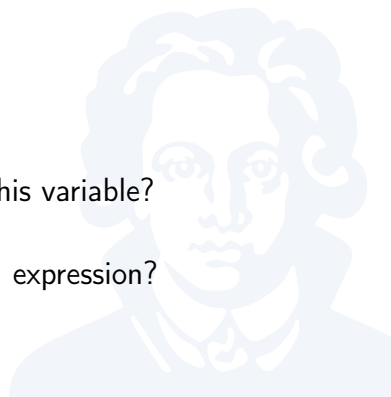


Dataflow Analysis

Analyse properties of variables or basic blocks

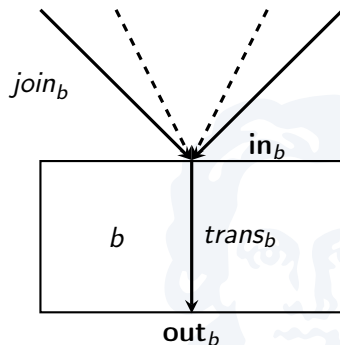
Examples in practice:

- Live Variables
Is this variable ever read?
- Reaching Definitions (def-use)
What are the possible values for this variable?
- Available Expressions
What variable definitely has which expression?

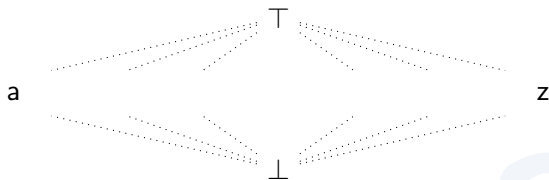


Formalisation

- \mathbf{in}_b : knowledge at entrance of basic block b
- \mathbf{out}_b : knowledge at exit of basic block b
- join_b : merges all \mathbf{out}_{b_i} for all basic blocks b_i that flow into b
- trans_b : updates \mathbf{out}_b from \mathbf{in}_b



Analysis lattice



- Analysis takes place in a *complete (semi)lattice*
- \top : 'Too much information / could be anything'
- \perp : 'No information'
- Join functions are typically join or meet operations
- Transfer functions $trans_b$ are *monotonic* (order-preserving)

if $a \leq c$, then $t(a) \leq trans_b(c)$

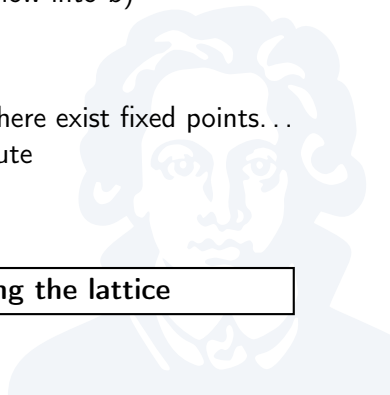
Computing the Least Fixed Point

Computing the results:

- Iterate updating
 - $\mathbf{in}_b := \text{join}_b(\mathbf{in}_{b_1}, \dots, \mathbf{in}_{b_n})$ (b_i flow into b)
 - $\mathbf{out}_b := \text{trans}_b(\mathbf{in}_b)$
- Until nothing changes

Knaster-Tarski theorem guarantees that there exist fixed points. . .
. . .but these may be impractical to compute

Take care when constructing the lattice



Example: Reaching Definitions

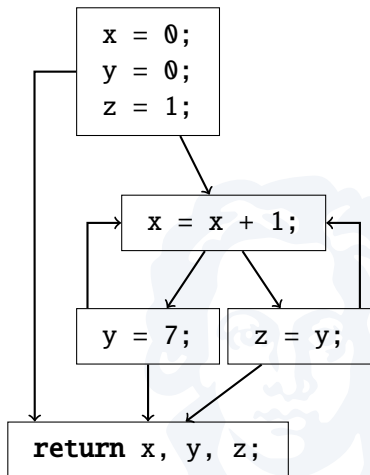
```

x = 0;
y = 0;
z = 1;

while (x < 5) {
  x = x + 1
  if (x >= 2) {
    y = 7;
  } else {
    z = y;
  }
}

return x, y, z;

```

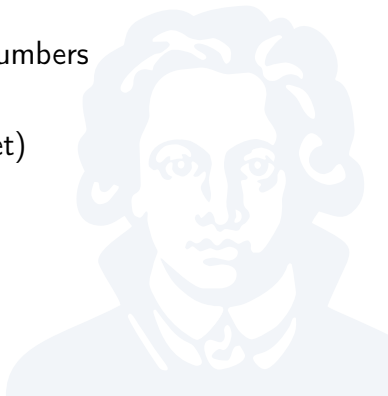


Reaching Definitions: What values are possible?

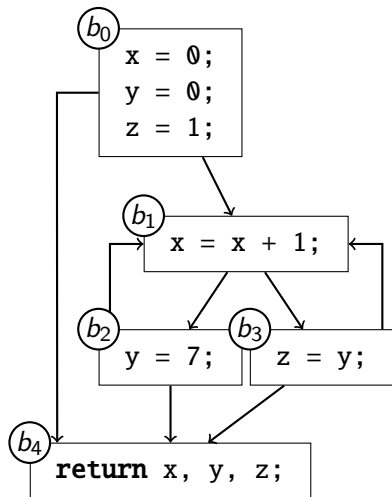
Example: Reaching Definitions

Designing our lattice:

- Capture sets of up to 3 possible numbers
- \top : More than 3 possible numbers
- \perp : \emptyset (no possible numbers seen yet)



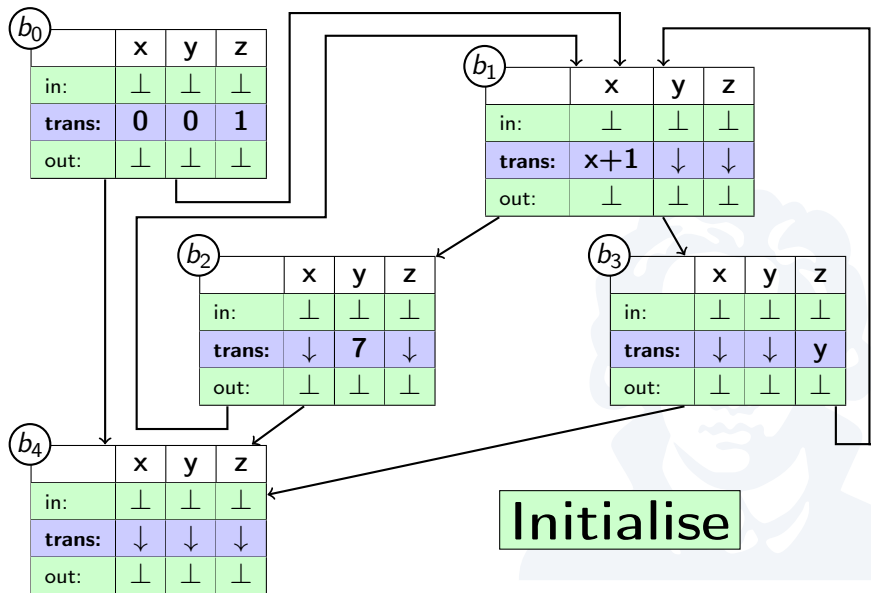
Example: Control-Flow Graph



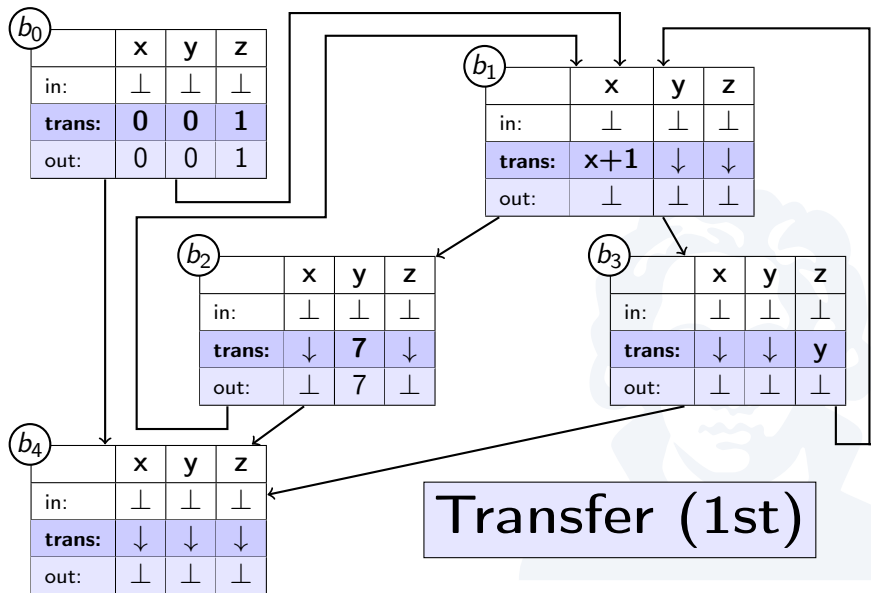
b	inputs	$trans_b$		
		x	y	z
b_0	\emptyset	0	0	1
b_1	$\{b_0, b_2, b_3\}$	$x + 1$	y	z
b_2	$\{b_1\}$	x	7	z
b_3	$\{b_1\}$	x	y	y
b_4	$\{b_0, b_2, b_3\}$	x	y	z

$$join_b = \bigcup_{s \in inputs_b} s$$

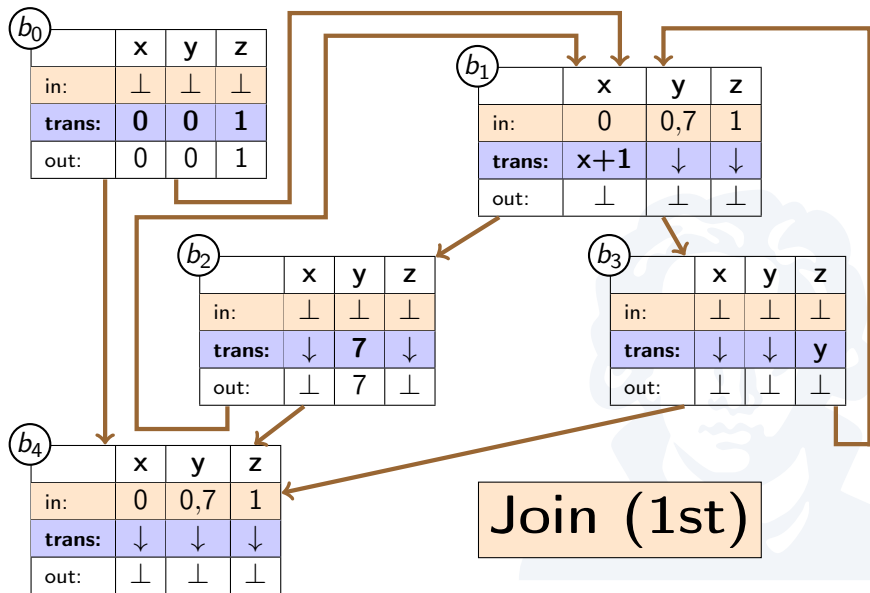
Example: Computing the Fixpoint



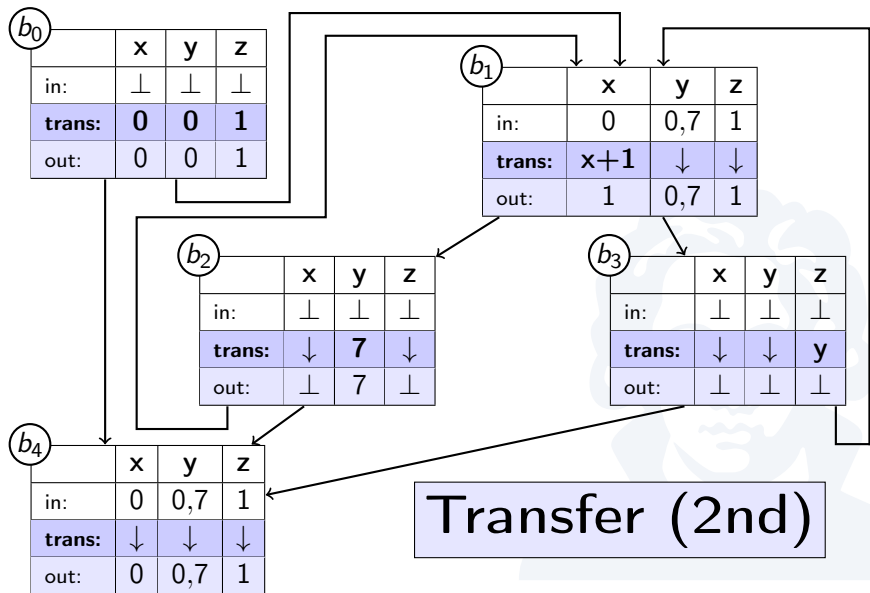
Example: Computing the Fixpoint



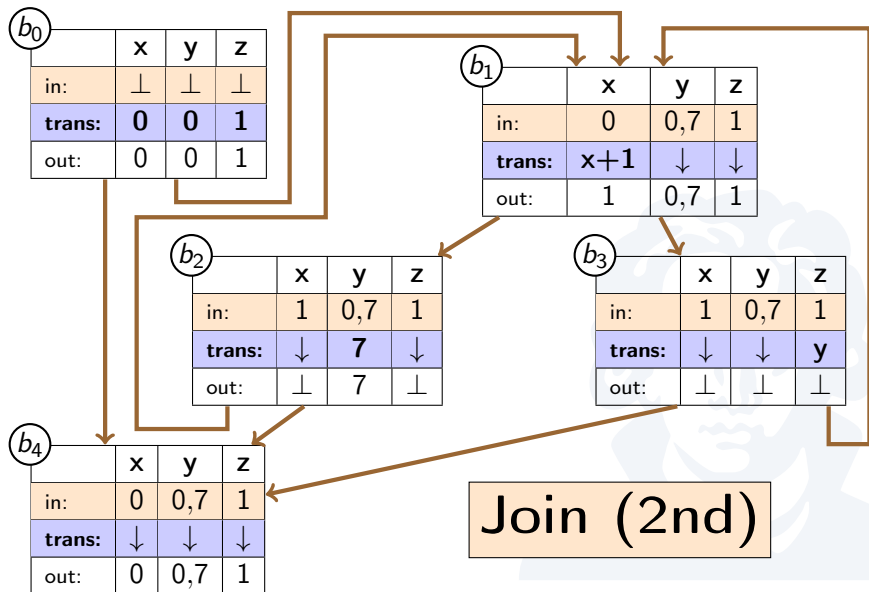
Example: Computing the Fixpoint



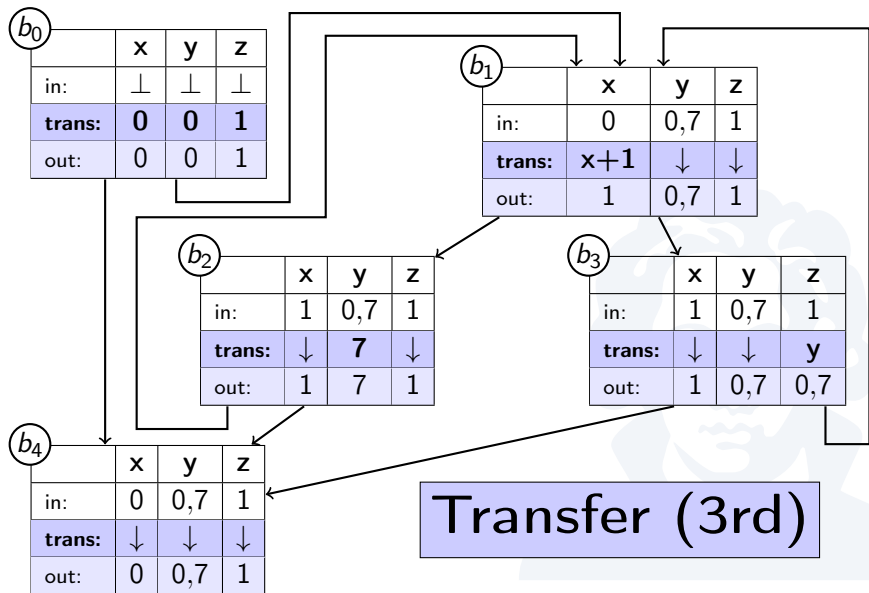
Example: Computing the Fixpoint



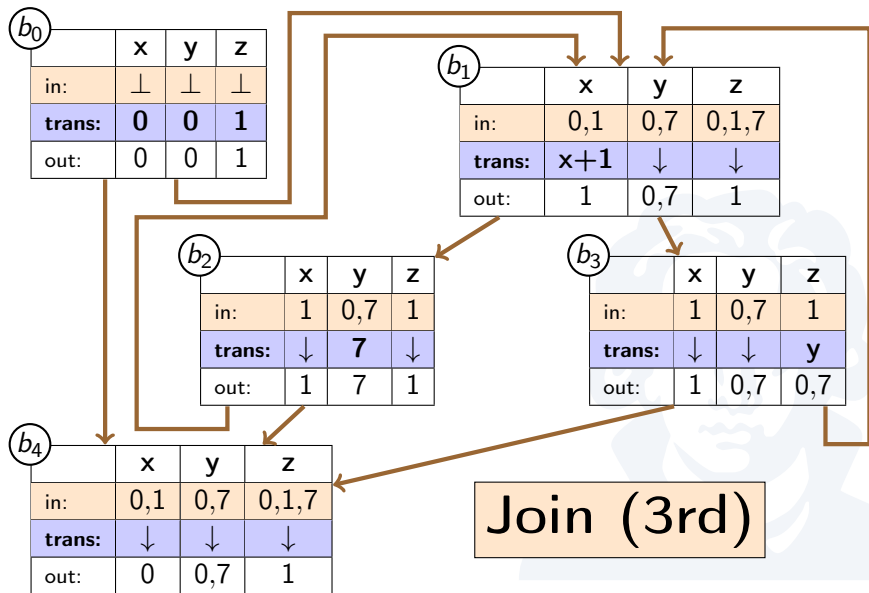
Example: Computing the Fixpoint



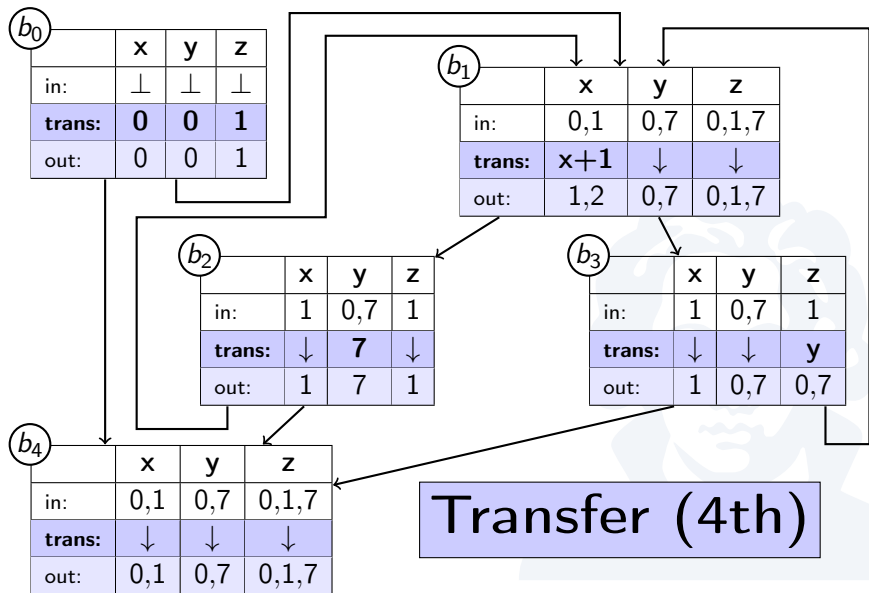
Example: Computing the Fixpoint



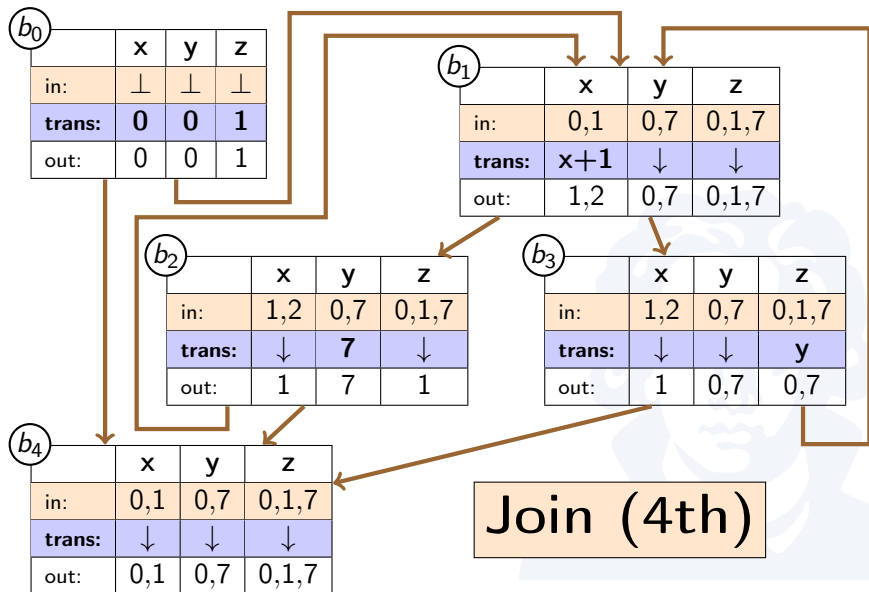
Example: Computing the Fixpoint



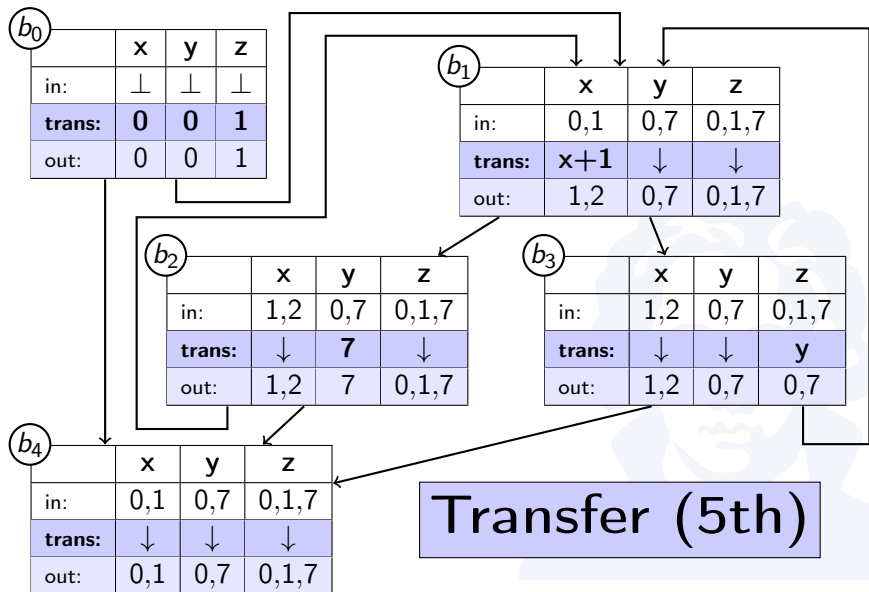
Example: Computing the Fixpoint



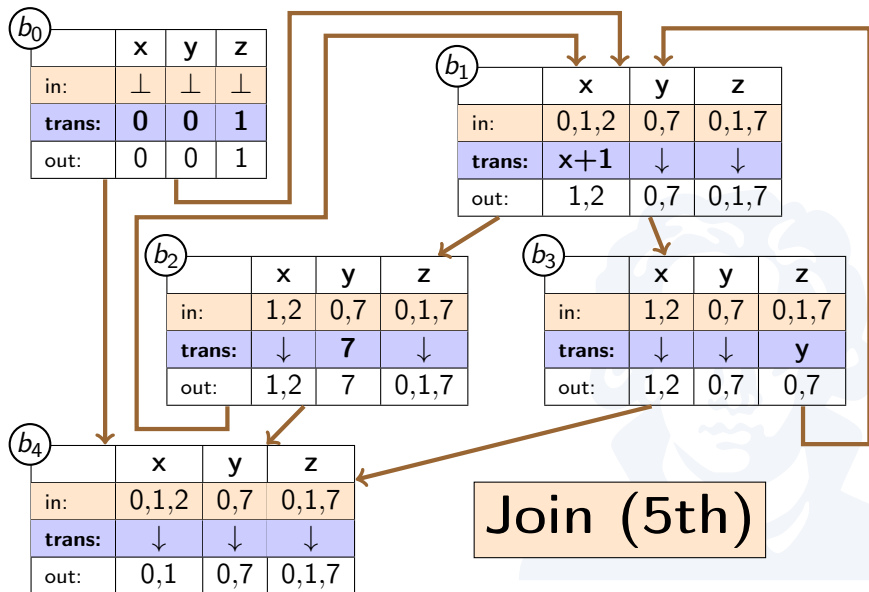
Example: Computing the Fixpoint



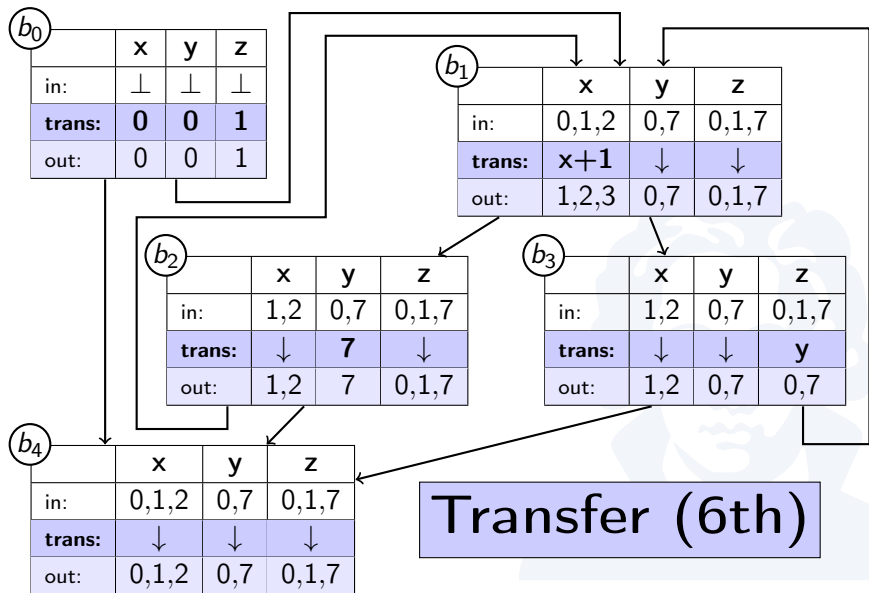
Example: Computing the Fixpoint



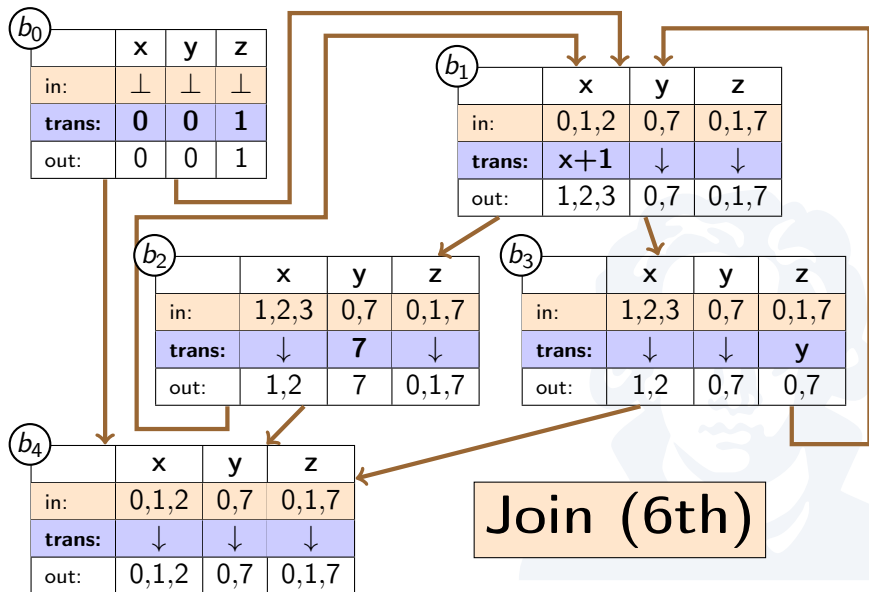
Example: Computing the Fixpoint



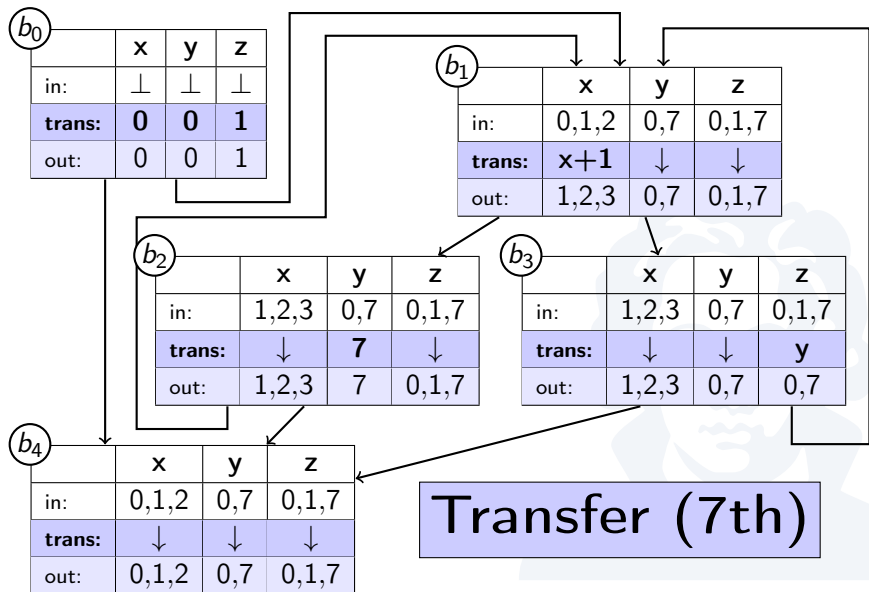
Example: Computing the Fixpoint



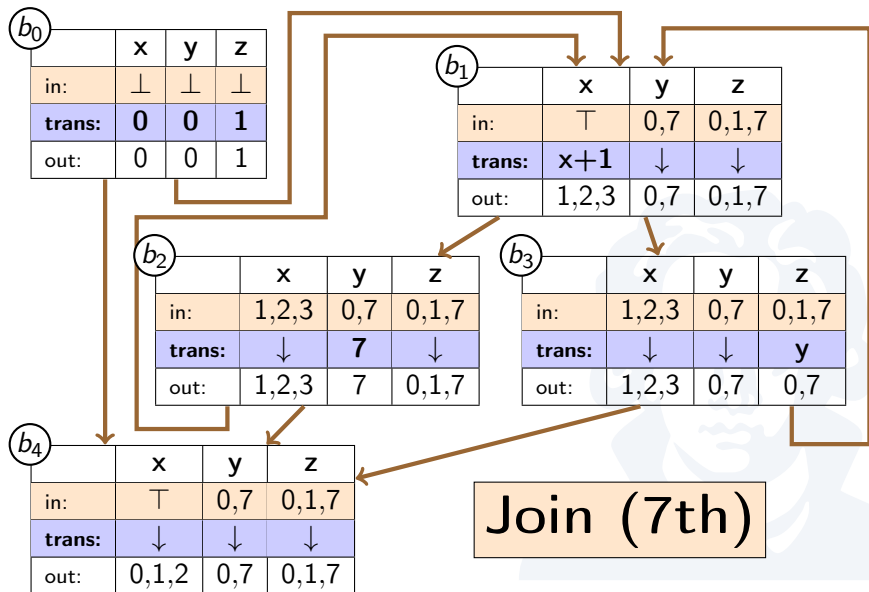
Example: Computing the Fixpoint



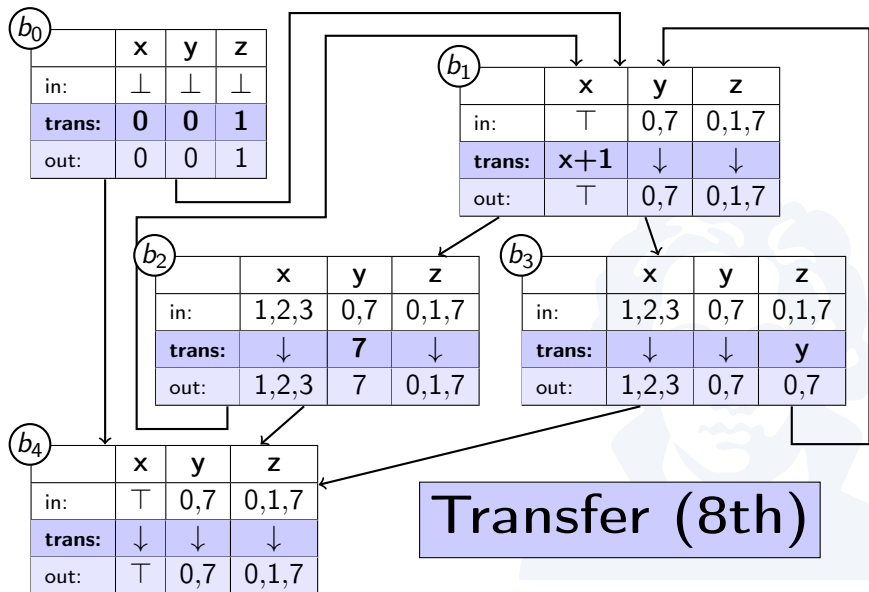
Example: Computing the Fixpoint



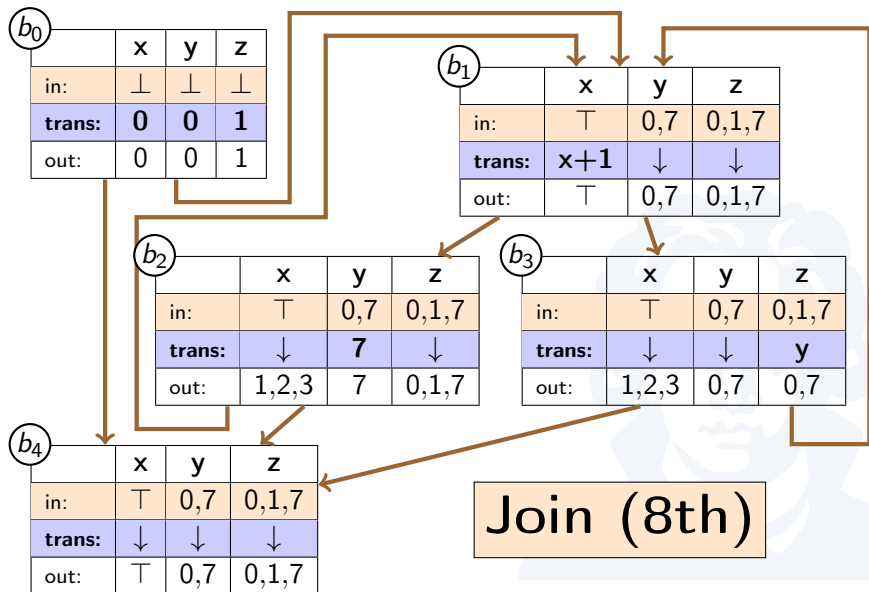
Example: Computing the Fixpoint



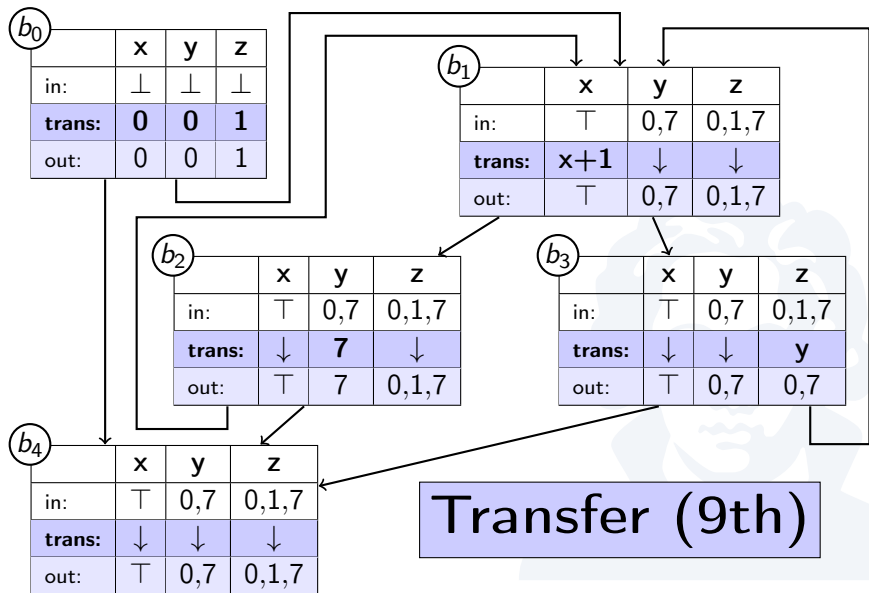
Example: Computing the Fixpoint



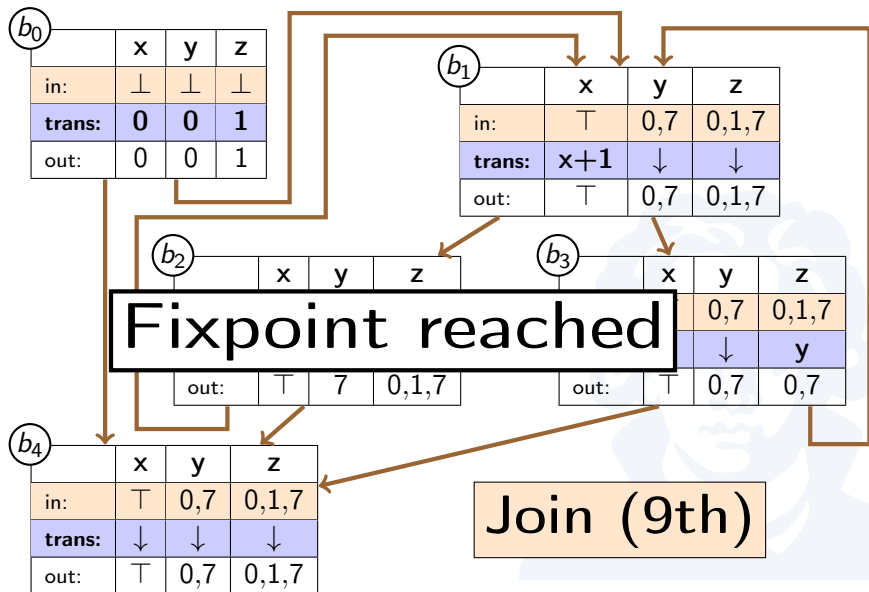
Example: Computing the Fixpoint



Example: Computing the Fixpoint



Example: Computing the Fixpoint



Example: Conclusion

```

x = 0;
y = 0;
z = 1;

while (x < 5) {
  x = x + 1
  if (x >= 2) {
    y = 7;
  } else {
    z = y;
  }
}

return x, y, z;

```

- Reached fixpoint after 9 iterations
- Return values:
 - x : \top (unknown/any)
 - y : 0 or 7
 - z : 0 or 1 or 7
- Conservative approximation of reality
- Once x reached more than 3 values, algorithm gave up and went to \top

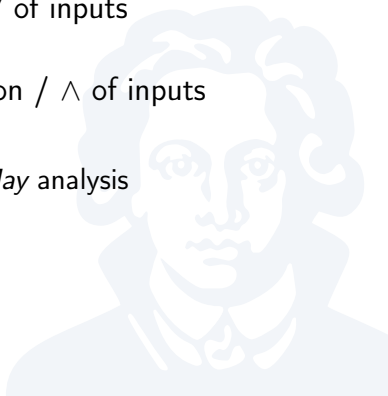
Alternatives are possible:

- Faster $x + 1$: always \top
- More precise conditions (flow-sensitive analysis)
- Bigger/smaller domain

“May” and “Must” analyses

- *May analysis*: Construct union / \vee of inputs
More information \rightarrow closer to \top
- *May analysis*: Construct intersection / \wedge of inputs
More information \rightarrow closer to \perp

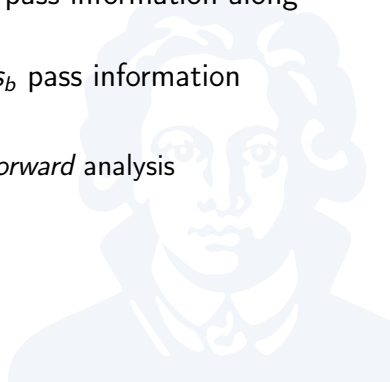
Reaching definitions is an example of a *May* analysis



Forward and Backward analyses

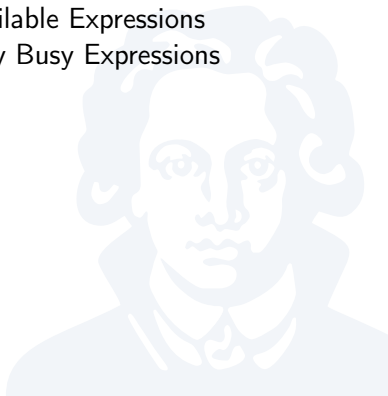
- *Forward analysis*: $join_b$ and $trans_b$ pass information along control flow graph
- *Backward analysis*: $join_b$ and $trans_b$ pass information along *inverse* control flow graph

Reaching definitions is an example of a *Forward* analysis



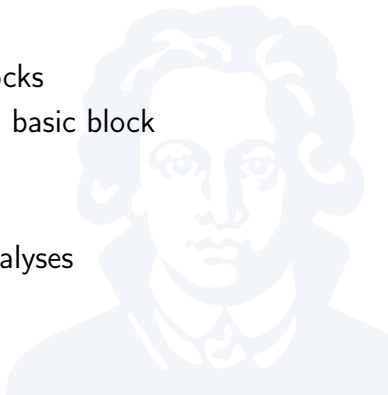
	May	
Forward		Reaching Definitions
Backward		Live Variables

	Must	
		Available Expressions
		Very Busy Expressions



Summary: Data-Flow Analysis

- Analyses properties of variables
- Analysis lattice encodes the viable values we consider
- Operates on control-flow graph
- $join_b$ merges inputs from input blocks
- $trans_b$ transfers information within basic block
- Compute fixed point
- Distinction: May/Must analyses
- Distinction: Forward/Backward analyses



Constraint-Based Analysis

- Idea: separate
 - *constraint generation* and
 - *constraint solution*
- Very generic strategy
- Easy to scale up and down
- Very little general structure



Abstract Interpretation

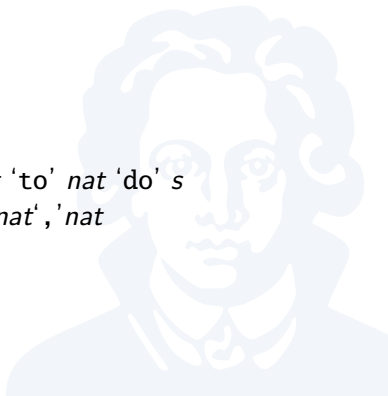
- Idea: use *alternative denotational semantics*, $\llbracket - \rrbracket_{\mathcal{A}}$
- $\llbracket - \rrbracket_{\mathcal{A}}$ is an *abstraction* over $\llbracket - \rrbracket$
- $\llbracket - \rrbracket_{\mathcal{A}}$ works on *abstract domain* that *simulates* real semantics

Consider the following language:

```

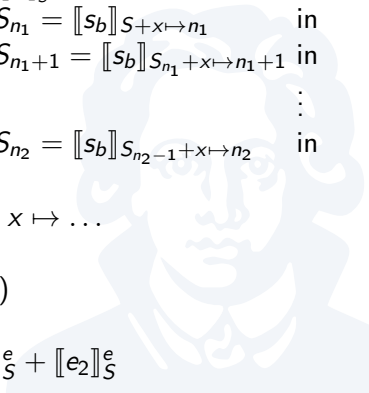
s ::= name ':=' e
   | 'begin' <s> ';' <s> 'end'
   | 'for' name 'from' nat 'to' nat 'do' s
   | 'read' name 'between' nat ',' nat

e ::= name
   | nat
   | <e> '+' <e>
  
```



Denotational Semantics

$\llbracket - \rrbracket^e$ for expressions (\rightarrow values), $\llbracket - \rrbracket^s$ for statements (\rightarrow states).

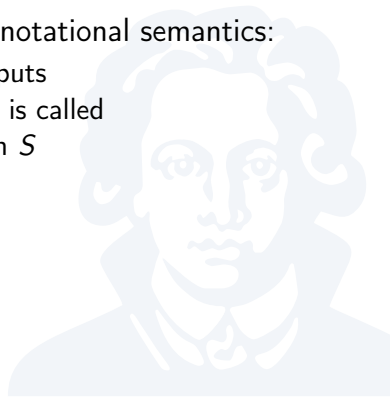
$$\begin{aligned}
 \llbracket x := e \rrbracket_S^s &= S + x \mapsto \llbracket e \rrbracket_S^e \\
 \llbracket \text{begin } s_1; s_2 \text{ end} \rrbracket_S^s &= \llbracket s_2 \rrbracket_{\llbracket s_1 \rrbracket_S^s}^s \\
 \llbracket \text{for } x \text{ from } n_1 \text{ to } n_2 \text{ do } s_b \rrbracket_S^s &= \text{let } S_{n_1} = \llbracket s_b \rrbracket_{S+x \mapsto n_1}^s \text{ in} \\
 &\quad \text{let } S_{n_1+1} = \llbracket s_b \rrbracket_{S_{n_1}+x \mapsto n_1+1}^s \text{ in} \\
 &\quad \vdots \\
 &\quad \text{let } S_{n_2} = \llbracket s_b \rrbracket_{S_{n_2-1}+x \mapsto n_2}^s \text{ in} \\
 &\quad S_{n_2} \\
 \llbracket \text{read } x \text{ between } n_1, n_2 \rrbracket_S^s &= S + x \mapsto \dots \\
 \llbracket x \rrbracket_S^e &= S(e) \\
 \llbracket n \rrbracket_S^e &= n \\
 \llbracket e_1 + e_2 \rrbracket_S^e &= \llbracket e_1 \rrbracket_S^e + \llbracket e_2 \rrbracket_S^e
 \end{aligned}$$


Denotational Semantics of inputs

Quick aside:

- To model inputs (as in **read** in denotational semantics):
 - Pass in (abstract) list of user inputs
 - 'Consume' input whenever **read** is called
 - Thread through program as with S

Omitted for brevity.



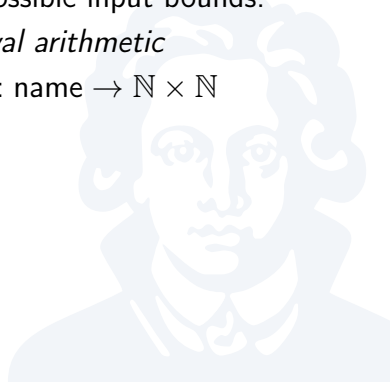
Abstract Domain Semantics: An Example

How can we make statements about such a program without having *concrete* user inputs?

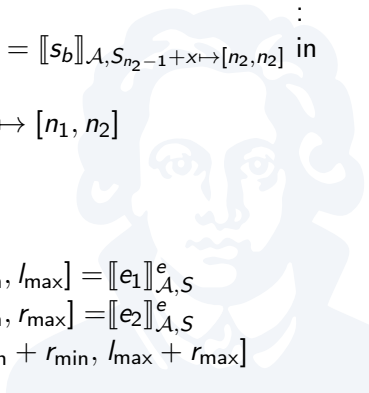
- Our **read** operation tells us the possible input bounds:
- Abstract interpretation over *interval arithmetic*
- For states S for $\llbracket - \rrbracket^e$ we have: $S : \text{name} \rightarrow \mathbb{N} \times \mathbb{N}$
- Notation: $[\text{min}, \text{max}]$

Example: $x := 2$ gives

$$S(x) = [2, 2]$$

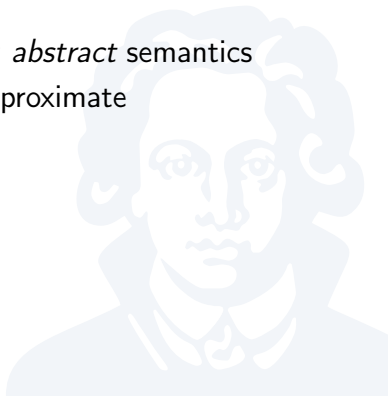


Formalising interval arithmetic

$$\begin{aligned}
 \llbracket x := e \rrbracket_{\mathcal{A}, S}^s &= S + x \mapsto \llbracket e \rrbracket_{\mathcal{A}, S}^e \\
 \llbracket \text{begin } s_1; s_2 \text{ end} \rrbracket_{\mathcal{A}, S}^s &= \llbracket s_2 \rrbracket_{\mathcal{A}, \llbracket s_1 \rrbracket_{\mathcal{A}, S}^s}^s \\
 \llbracket \text{for } x \text{ from } n_1 \text{ to } n_2 \text{ do } s_b \rrbracket_{\mathcal{A}, S}^s &= \text{let } S_{n_1} = \llbracket s_b \rrbracket_{\mathcal{A}, S+x \mapsto [n_1, n_1]}^s \quad \text{in} \\
 &\quad \vdots \\
 &\quad \text{let } S_{n_2} = \llbracket s_b \rrbracket_{\mathcal{A}, S_{n_2-1}+x \mapsto [n_2, n_2]}^s \quad \text{in} \\
 &\quad S_{n_2} \\
 \llbracket \text{read } x \text{ between } n_1, n_2 \rrbracket_{\mathcal{A}, S}^s &= S + x \mapsto [n_1, n_2] \\
 \llbracket x \rrbracket_{\mathcal{A}, S}^e &= S(e) \\
 \llbracket n \rrbracket_{\mathcal{A}, S}^e &= [n, n] \\
 \llbracket e_1 + e_2 \rrbracket_{\mathcal{A}, S}^e &= \text{let } [l_{\min}, l_{\max}] = \llbracket e_1 \rrbracket_{\mathcal{A}, S}^e \\
 &\quad [r_{\min}, r_{\max}] = \llbracket e_2 \rrbracket_{\mathcal{A}, S}^e \\
 &\quad \text{in } [l_{\min} + r_{\min}, l_{\max} + r_{\max}]
 \end{aligned}$$


Summary: Abstract Interpretation

- Replace denotational semantics by *abstract* semantics
- Abstract semantics should over-approximate



Analysis Qualities: Soundness

Soundness:

- Does the analysis correctly model *all* possibilities?
- Unsound \Rightarrow false negatives
- Compiler optimisations *must* be fully sound
- Unsoundness is permissible for:
 - *Static checkers*: May miss some bugs
 - *Speculative optimisations*:
 - Optimised code has *guard*, validates claims at run-time
 - Revert to unoptimised behaviour otherwise

Unsoundness example:

```
int i = 2;  
if (...) { i = 3; }
```

Claim: "i is always 3"

Analysis Qualities: Precision

Precision:

- Does the analysis model *only real* possibilities?
- Imprecise \Rightarrow false positives
- Imprecision may forbid some viable optimisations

Conservative analysis: sound but not 100% precise

Imprecision example:

```
int i = 2;  
if (...) { i = 3; }
```

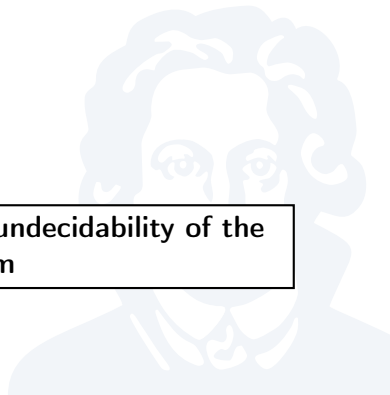
Claim: “i is between 0 and 1000”

Analysis Trade-offs

Analyses are trade-offs between:

- Soundness
- Precision
- Execution time

Soundness/precision bounded by undecidability of the Halting Problem



Analysis Qualities

- Intraprocedural Analyses
Only within individual procedures
- Interprocedural Analyses
Across procedure boundaries

Properties of particular analyses:

- Context Sensitivity
- Flow Sensitivity
- Path Sensitivity



Context Sensitivity

Consider *Calling context*:

```
int f(int i)
{
  return i;
}

...
x = 2;
y = f(x);
```

With context sensitivity

- We have one `i` per call site of `f`
- `y` is 2

Without context sensitivity

- We have one `i` total
- `y` is some `int`

Analogous to *inlining*

Flow Sensitivity

Consider *Control flow* and *Order of execution*:

With flow sensitivity

- y is 2

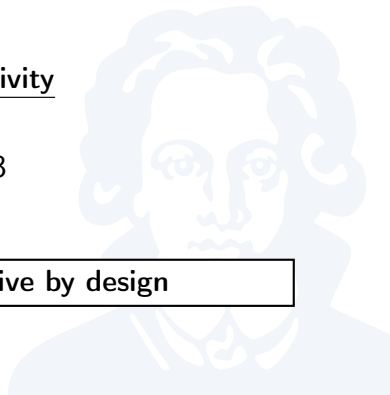
Without flow sensitivity

(assuming no SSA)

- y may be 2 or 3

```
x = 2;  
y = x;  
x = 3;
```

Flow analysis is flow sensitive by design



Path Sensitivity

Consider properties inferred from *Execution path*:

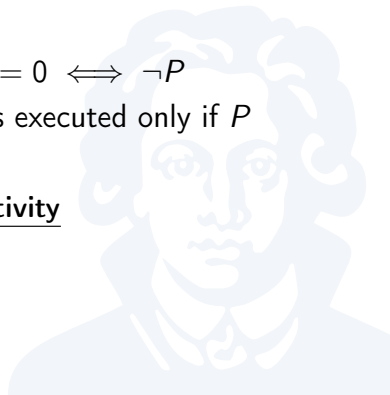
With path sensitivity

- $y \in \{1, 2\}$
- Records that $x = 0 \iff \neg P$
- Knows that A is executed only if P
 $\Rightarrow x \neq 0$ at A

Without path sensitivity

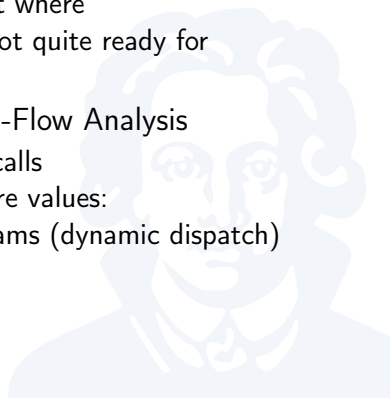
- $y \in \{0, 1, 2\}$
- Less precise

```
x = 0;
if (P) {
  x = 1;
}
y = 2;
if (P) {
  y = x; /*A*/
}
```



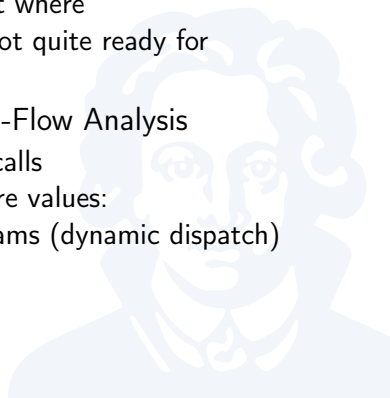
Other important analyses

- Points-to analysis
 - Determines which pointers point where
 - Substantial research area, still not quite ready for practical use
- k-CFA: (context sensitive) Control-Flow Analysis
 - Determines targets of function calls
 - Important whenever functions are values:
Functional programs, OO programs (dynamic dispatch)



Other important analyses

- Points-to analysis
 - Determines which pointers point where
 - Substantial research area, still not quite ready for practical use
- k-CFA: (context sensitive) Control-Flow Analysis
 - Determines targets of function calls
 - Important whenever functions are values:
Functional programs, OO programs (dynamic dispatch)



Summary

- Analysis qualities
- Precision & Soundness
- Flow-, Path-, Context-sensitivity

Kinds of analysis:

- Type and Effect analysis
- Flow Analysis
- Constraint-Based Analysis
- Abstract Interpretation



Literature

- Flemming Nielson, Hanne R. Nielson, Chris Hankin: “Principles of Program Analysis” (book)
- Alfred V. Aho, Monica S. Lam, Ravi Sethi, Jeffrey D. Ullman: “Compilers: Principles, Techniques, and Tools” (book)
- Alex Aiken: “Introduction to Set Constraint-Based Program Analysis”
- Jeanne Ferrante, Karl J. Ottenstein, Joe D. Warren: “The Program Dependence Graph and Its Use in Optimization” (ACM TOPLAS, 1987)
- Sumit Gulwani, Ashish Tiwari: “Computing Procedure Summaries for Interprocedural Analysis ” (ESOP 2007)
- Ondej Lhoták , Yannis Smaragdakis, and Manu Sridharan: “Pointer Analysis”, (Report from Dagstuhl Seminar 13162)

Next week:

Dynamic Program Analysis

