

LESSON 2

DATA FLOW-BASED ANALYSIS

STATIC PROGRAM ANALYSIS AND CONSTRAINT SOLVING

MASTER'S DEGREE IN FORMAL METHODS IN COMPUTER SCIENCE

Manuel Montenegro (montenegro@fdi.ucm.es)

Departamento de Sistemas Informáticos y Computación

Facultad de Informática

Universidad Complutense de Madrid

1. Introduction to static analysis
2. Data-flow analysis
3. Live variable analysis
4. Available expressions analysis
5. Monotone frameworks

- Introduction to lattices

- Lattice construction

- Instances of monotone frameworks

INTRODUCTION TO STATIC ANALYSIS

- It addresses the properties a program may satisfy.
 - Does a function return an integer value?
 - Does a program terminate for any input?
 - Would a program try to access dangling pointers?
 - Would variables x and y refer to the same memory location?
- **Static analysis:** it is carried out without executing the program.

In which ways do analysis behave?

*Does P hold
for the input program?*

*Does the analysis
report P ?*

	YES	NO
YES	✓	False positive
NO	False negative	✓

*Does P hold
for the input program?*

*Does the analysis
report P ?*

	YES	NO
YES	✓	
NO		✓

Theorem

*Every non-trivial and extensional property is **undecidable**.*

- **Non-trivial:** Some programs satisfy the property, some others do not.
- **Extensional:** Replacing a program with an equivalent one preserves the property.

*Does P hold
for the input program?*

*Does the analysis
report P ?*

	YES	NO
YES	✓	
NO		✓

- If we want to achieve this, we have to relax our expectations:
 - Analysis with manual intervention, or
 - Analysis constrained to a decidable subset of the target language, or
 - Analysis of programs that finish after a finite amount of computation steps.

WHAT WE COULD ACHIEVE: UPPER APPROXIMATION

*Does P hold
for the input program?*

*Does the analysis
report P ?*

	YES	NO
YES	✓	False positive
NO		✓

- The analysis returns either *No* or *Don't know*.
- Used to detect **unpleasant** properties in a program: dangling pointers, type-related inconsistencies, unexpected sharing, etc.

WHAT WE COULD ACHIEVE: LOWER APPROXIMATION

*Does P hold
for the input program?*

*Does the analysis
report P ?*

	YES	NO
YES	✓	
NO	False negative	✓

- The analysis returns either *Yes* or *Don't know*.
- Used to detect **desirable** properties in programs, e.g., termination.

- **Data-flow analysis**
 - Gather values computed in every execution point.
- **Abstract interpretation**
 - Approximate execution of programs by using properties as computable values.
- **Type systems**
 - Classify computable values into types, and check for discrepancies.
- **Symbolic execution**
 - Compute those inputs that make the program take a given path.
- **Constraint-based analysis**
- **Model checking**
- **Hoare's logic**, etc.

DATA-FLOW ANALYSIS

- A **data-flow analysis** obtains information on the values computed by a program in each execution point.
- It receives the *Control-Flow Graph* (CFG) of the program being analysed.
- More info on CFGs and how to build them:



S. Muchnick

Advanced Compiler Design and Implementation

Morgan-Kaufmann (1998)

THE WHILE LANGUAGE

- Variables: $x, y, z \in \text{Var}$.
- Numbers: $m \in \text{Num}$.
- Arithmetic expressions (**Exp**):

$$e ::= x \mid m \mid e_1 + e_2 \mid e_1 * e_2 \mid \dots$$

- Boolean expressions:

$$b ::= \text{true} \mid \text{false} \mid \text{not } b \mid b_1 \text{ and } b_2 \mid \dots$$

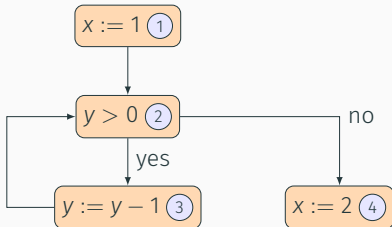
- Statements (**Stm**):

$$\begin{aligned} S &::= \text{skip} \\ &\mid x := e \\ &\mid S_1; S_2 \\ &\mid \text{if } b \text{ then } S_1 \text{ else } S_2 \\ &\mid \text{while } b \text{ do } S \end{aligned}$$

CONTROL-FLOW GRAPHS

- We assume that programs have been transformed into a CFG, where the basic blocks are **skip** statements, assignments, and boolean expressions.
- Each basic block will contain a unique **label** taken from the set **Lab** = {1, 2, 3, ...}.

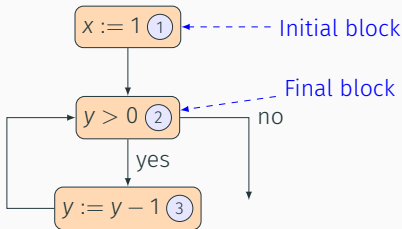
```
x := 1;  
  
while y > 0 do  
    y := y - 1;  
  
x := 2;
```



CONTROL FLOW GRAPHS

- Each program has an **initial block** and a **final block**.
- We assume that the initial block has no predecessors.
 - If this does not hold, add a **skip** block at the beginning.

```
x := 1;  
  
while y > 0 do  
    y := y - 1;
```



- **Context sensitive/insensitive.**
 - Depending on whether it takes into account the call stack leading to a given program point.
- **Path sensitive/insensitive.**
 - Depending on whether identifies the conditions that lead the execution to take a given path.
- **Control-flow sensitive/insensitive.**
 - Depending if they consider the order of the basic statements in the program.

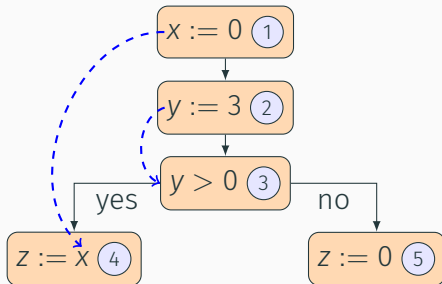
- **Must analysis**
 - It determines whether a property holds in **every** execution path of the program.
- **May analysis**
 - It determines whether a property holds in **some** execution path of the program.

- **Forward analysis**
 - The information on a program point is computed from the program points that have been executed before.
- **Backward analysis**
 - The information on a program point is computed from the program points that might be reached later.

LIVE VARIABLE ANALYSIS

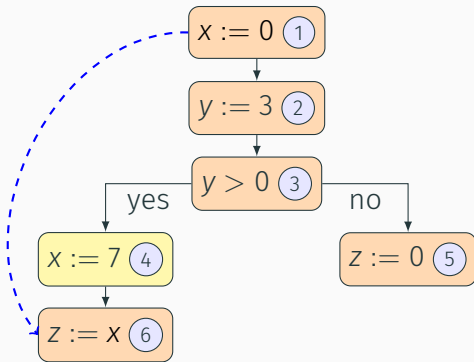
LIVE VARIABLE ANALYSIS

- A variable is **live** at some point n if there is an execution path from n to the end of the program in which x is read before being assigned another value.
 - In other words, x is live at some point if its value at that point may be needed in the future.



- x is live at the end of block 1.
- y is live at the end of block 2.
- z is not live in any block.

LIVE VARIABLE ANALYSIS



- `x` is **not** live at the end of block 1.
- ...because its occurrence in block 5 is preceded by an assignment to `x` in block 4.

WHY IS IT USEFUL?

- **Register allocation**

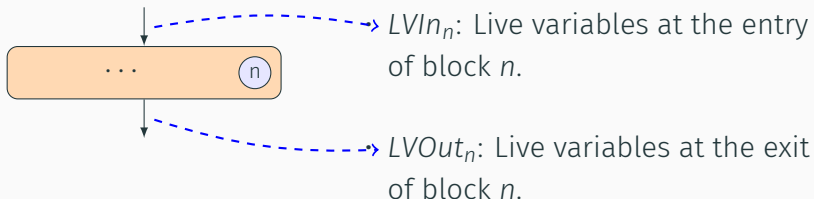
If two variables are never live at the same point, they can be stored in the same machine register.

- **Dead code elimination**

Any assignment to a variable that is not live at the end of the block can be removed.

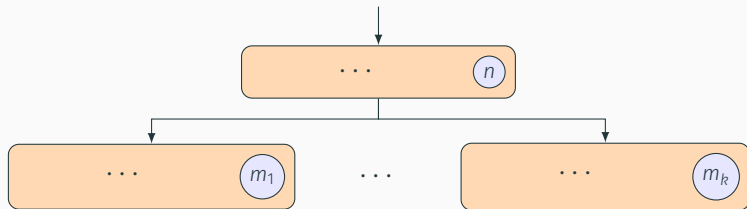
LIVE VARIABLE ANALYSIS: BACKWARD TRAVERSAL

- **Aim:** Determine which variables **may** be live at each program point.
- It will be a **backward analysis**.
 - The information in a given block will be obtained from the blocks executed after it.
- We compute, for each block:



LIVE VARIABLE ANALYSIS: DATA-FLOW EQUATIONS

- A variable is live at the **exit** of a block if it is live at the **entrance** of any of the blocks following it (successors).



$$LVOut_n = \bigcup_{m \in \text{succ}(n)} LVIn_m$$

- If n is a final block, then $LVOut_n = \emptyset$.

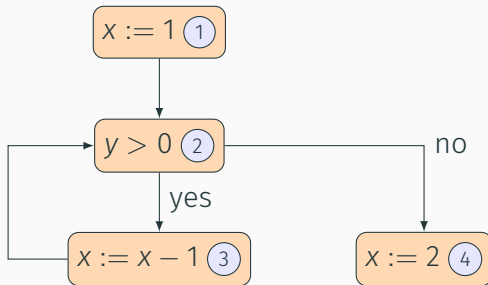
LIVE VARIABLE ANALYSIS: DATA-FLOW EQUATIONS

- A variable is live at the **entry** of a block if:
 - This variable is read in that block, or
 - It was live at the end of the block, and it has been assigned to during its execution.
- Let us define:
 - Gen_n : Variables that are read in block n .
 - $Kill_n$: Variables that are written to in block n .
 - If block n contains an assignment $x := e$, then $Kill_n = \{x\}$. Otherwise $Kill_n = \emptyset$.

Dataflow equation:

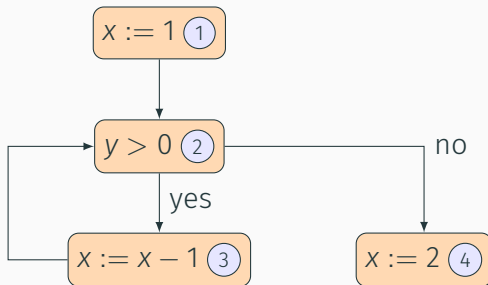
$$LVIn_n = (LVOut_n - Kill_n) \cup Gen_n$$

EXAMPLE



Gen_1	$= \emptyset$	$Kill_1$	$= \{x\}$
Gen_2	$= \{y\}$	$Kill_2$	$= \emptyset$
Gen_3	$= \{x\}$	$Kill_3$	$= \{x\}$
Gen_4	$= \emptyset$	$Kill_4$	$= \{x\}$

EXAMPLE



$$LVOut_1 = LVIn_2$$

$$LVOut_2 = LVIn_3 \cup LVIn_4$$

$$LVOut_3 = LVIn_2$$

$$LVOut_4 = \emptyset$$

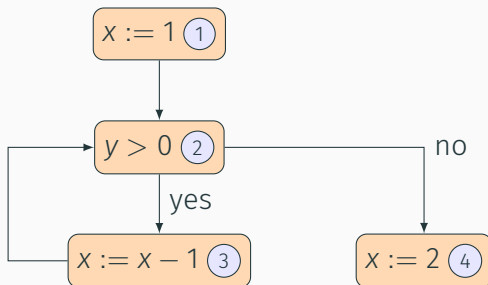
$$LVIn_1 = (LVOut_1 - \{x\}) \cup \emptyset$$

$$LVIn_2 = (LVOut_2 - \emptyset) \cup \{y\}$$

$$LVIn_3 = (LVOut_3 - \{x\}) \cup \{x\}$$

$$LVIn_4 = (LVOut_4 - \{x\}) \cup \emptyset$$

EXAMPLE



$$LVOut_1 = LVIn_2$$

$$LVOut_2 = LVIn_3 \cup LVIn_4$$

$$LVOut_3 = LVIn_2$$

$$LVOut_4 = \emptyset$$

$$LVIn_1 = LVOut_1 - \{x\}$$

$$LVIn_2 = LVOut_2 \cup \{y\}$$

$$LVIn_3 = LVOut_3 \cup \{x\}$$

$$LVIn_4 = LVOut_4 - \{x\}$$

EQUATION SYSTEMS

- The result is a **system of equations** whose unknown variables are $(LVOut_1, \dots, LVOut_4, LVIn_1, \dots, LVIn_4)$.

$$\begin{array}{ll} LVOut_1 = LVIn_2 & LVIn_1 = LVOut_1 - \{x\} \\ LVOut_2 = LVIn_3 \cup LVIn_4 & LVIn_2 = LVOut_2 \cup \{y\} \\ LVOut_3 = LVIn_2 & LVIn_3 = LVOut_3 \cup \{x\} \\ LVOut_4 = \emptyset & LVIn_4 = LVOut_4 - \{x\} \end{array}$$

- Solutions of this system are subsets of the set of variables in the program. In our case: **Var** = $\{x, y\}$.
- Here we have a solution:

$$\begin{array}{llll} LVOut_1 = \{x, y\} & LVOut_2 = \{x, y\} & LVOut_3 = \{x, y\} & LVOut_4 = \emptyset \\ LVIn_1 = \{y\} & LVIn_2 = \{x, y\} & LVIn_3 = \{x, y\} & LVIn_4 = \emptyset \end{array}$$

- Live variable analysis requires solving a system of equations which involves operations on set.
- Pending questions:
 1. Does the system have a solution?
 2. Could it have several solutions?
 3. In case it has several solutions, is there a solution that is better than the others?
 4. In that case, how can we compute the best solution?

SOLVING SYSTEMS OF EQUATIONS

- Let us define $F : \mathcal{P}(\text{Var})^8 \rightarrow \mathcal{P}(\text{Var})^8$ as follows:

$$F \begin{pmatrix} LVOut_1 \\ LVOut_2 \\ LVOut_3 \\ LVOut_4 \\ LVIn_1 \\ LVIn_2 \\ LVIn_3 \\ LVIn_4 \end{pmatrix} = \begin{pmatrix} LVIn_2 \\ LVIn_3 \cup LVIn_4 \\ LVIn_2 \\ \emptyset \\ LVOut_1 - \{x\} \\ LVOut_2 \cup \{y\} \\ LVOut_3 \cup \{x\} \\ LVOut_4 - \{x\} \end{pmatrix}$$

- If **LV** denotes the following vector

$$\mathbf{LV} = (LVOut_1, \dots, LVOut_4, LVIn_1, \dots, LVIn_4)^T$$

we can formalize the equation system as follows:

$$\mathbf{LV} = F(\mathbf{LV})$$

- Given a function $F : L \rightarrow L$, we say that $x \in L$ is a **fixed point** of F if $F(x) = x$.

$$LV = F(LV)$$

The solutions of this system are the fixed points of F

1. Does the system have a solution?
2. Could it have several solutions?
3. In case it has several solutions, is there a solution that is better than the others?
4. In that case, how can we compute the best solution?

- Given a function $F : L \rightarrow L$, we say that $x \in L$ is a **fixed point** of F if $F(x) = x$.

$$LV = F(LV)$$

The solutions of this system are the fixed points of F

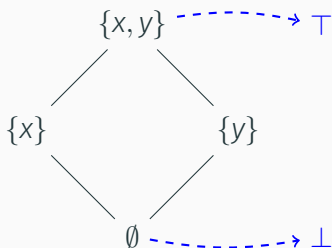
1. Does F have a fixed point?
2. Could it have several fixed points?
3. In case it has several fixed points, is there a fixed point that is better than the others?
4. In that case, how can we compute the best fixed point?

- In order to know whether a solution is better than the others, we have to set up an **order** between solutions.
- A **partial order** on a set L is a relation \sqsubseteq satisfying the following properties:
 - *Reflexive*: For all $x \in L$: $x \sqsubseteq x$.
 - *Antisymmetric*: For all $x, y \in L$: if $x \sqsubseteq y$, $y \sqsubseteq x$, then $x = y$.
 - *Transitive*: For all $x, y, z \in L$: if $x \sqsubseteq y$, $y \sqsubseteq z$, then $x \sqsubseteq z$.
- A **partially ordered set** (*poset*) is a set L together with an order relation \sqsubseteq on this set. It is denoted by (L, \sqsubseteq) .

Example

$(\mathcal{P}(\text{Var}), \subseteq)$ is an ordered set, where $\text{Var} = \{x, y\}$.

- Graphical representation of $(\mathcal{P}(\text{Var}), \subseteq)$:



- A poset is **bounded** if it has a pair of elements \top, \perp such that:
 - $\perp \subseteq x$ for all $x \in L$.
 - $x \subseteq \top$ for all $x \in L$.

TUPLES OF ORDERED SETS

- The solutions of our system are not elements of $\mathcal{P}(\mathbf{Var})$, but **tuples** of the set $\mathcal{P}(\mathbf{Var})^8$.
- If $(A_1, \sqsubseteq_1), (A_2, \sqsubseteq_2), \dots, (A_n, \sqsubseteq_n)$ are posets, we can define an order \sqsubseteq on $A_1 \times A_2 \times \dots \times A_n$ as follows:

$$(x_1, \dots, x_n) \sqsubseteq (y_1, \dots, y_n) \text{ iff } x_1 \sqsubseteq_1 y_1, x_2 \sqsubseteq_2 y_2, \dots, x_n \sqsubseteq_n y_n$$

Example

- In $\mathcal{P}(\{x, y\})^2$:

$$(\emptyset, \{x\}) \sqsubseteq (\{x\}, \{x\}) \quad (\emptyset, \emptyset) \sqsubseteq (\{x\}, \{y\})$$

but

$$(\emptyset, \{x\}) \not\sqsubseteq (\{y\}, \emptyset) \quad (\{y\}, \emptyset) \not\sqsubseteq (\emptyset, \{x\})$$

- Therefore, $(\mathcal{P}(\mathbf{Var})^8, \sqsubseteq)$ is a poset.
- Moreover, it is bounded. Its lower bound is:

$$\perp = (\underbrace{\emptyset, \emptyset, \dots, \emptyset}_{8 \text{ times}})$$

and its upper bound is:

$$\top = (\underbrace{\mathbf{Var}, \mathbf{Var}, \dots, \mathbf{Var}}_{8 \text{ times}})$$

MONOTONIC FUNCTIONS

- Given a poset (L, \sqsubseteq) , a function $F : L \rightarrow L$ is **monotonically increasing** if and only if for all $x, y \in L$:

$$x \sqsubseteq y \implies F(x) \sqsubseteq F(y)$$

- In our example, if F is defined as follows:

$$F \begin{pmatrix} LVOut_1 \\ LVOut_2 \\ LVOut_3 \\ LVOut_4 \\ LVIn_1 \\ LVIn_2 \\ LVIn_3 \\ LVIn_4 \end{pmatrix} = \begin{pmatrix} LVIn_2 \\ LVIn_3 \cup LVIn_4 \\ LVIn_2 \\ \emptyset \\ LVOut_1 - \{x\} \\ LVOut_2 \cup \{y\} \\ LVOut_3 \cup \{x\} \\ LVOut_4 - \{x\} \end{pmatrix}$$

Then it is monotonically increasing.

ASCENDING KLEENE CHAIN

- Let us start with an ordered set (L, \sqsubseteq) with a lower bound \perp and a monotonically increasing function $F : L \rightarrow L$.
- Obviously, $\perp \sqsubseteq F(\perp)$, by definition of \perp .
- Moreover, since F is monotonic, $F(\perp) \sqsubseteq F(F(\perp))$.
- Moreover, since F is monotonic, $F(F(\perp)) \sqsubseteq F(F(F(\perp)))$.
- By successively applying F to the \perp element, we obtain an **ascending chain**.

$$\perp \sqsubseteq F(\perp) \sqsubseteq F(F(\perp)) \sqsubseteq F(F(F(\perp))) \sqsubseteq \dots$$

- or, equivalently

$$\perp \sqsubseteq F(\perp) \sqsubseteq F^2(\perp) \sqsubseteq F^3(\perp) \sqsubseteq \dots \sqsubseteq F^n(\perp) \sqsubseteq \dots$$

ASCENDING KLEENE CHAIN

$$F \begin{pmatrix} LVOut_1 \\ LVOut_2 \\ LVOut_3 \\ LVOut_4 \\ LVIIn_1 \\ LVIIn_2 \\ LVIIn_3 \\ LVIIn_4 \end{pmatrix} = \begin{pmatrix} LVIIn_2 \\ LVIIn_3 \cup LVIIn_4 \\ LVIIn_2 \\ \emptyset \\ LVOut_1 - \{x\} \\ LVOut_2 \cup \{y\} \\ LVOut_3 \cup \{x\} \\ LVOut_4 - \{x\} \end{pmatrix}$$

We obtain the following chain:

$$\underbrace{\begin{pmatrix} \emptyset \\ \emptyset \\ \emptyset \\ \emptyset \\ \emptyset \\ \emptyset \\ \emptyset \\ \emptyset \end{pmatrix}}_{\perp} \sqsubseteq \underbrace{\begin{pmatrix} \emptyset \\ \emptyset \\ \emptyset \\ \emptyset \\ \emptyset \\ \{y\} \\ \{x\} \\ \emptyset \end{pmatrix}}_{F(\perp)} \sqsubseteq \underbrace{\begin{pmatrix} \{y\} \\ \{x\} \\ \{y\} \\ \emptyset \\ \emptyset \\ \{y\} \\ \{x\} \\ \emptyset \end{pmatrix}}_{F^2(\perp)} \sqsubseteq \underbrace{\begin{pmatrix} \{y\} \\ \{x\} \\ \{y\} \\ \emptyset \\ \{y\} \\ \{x, y\} \\ \{x, y\} \\ \emptyset \end{pmatrix}}_{F^3(\perp)} \sqsubseteq \underbrace{\begin{pmatrix} \{x, y\} \\ \{x, y\} \\ \{x, y\} \\ \emptyset \\ \{y\} \\ \{x, y\} \\ \{x, y\} \\ \emptyset \end{pmatrix}}_{F^4(\perp)} \sqsubseteq \dots$$



CAN THIS CHAIN INCREASE INDEFINITELY?

- Therefore, our chain has the following behaviour:

$$\perp \sqsubset F(\perp) \sqsubset F^2(\perp) \sqsubset \dots \sqsubset F^k(\perp) = F^{k+1}(\perp) = F^{k+2}(\perp) = \dots$$

- We say that the chain **stabilizes** after the k -th iteration.
- Since $F^k(\perp) = F^{k+1}(\perp) = F(F^k(\perp))$, we get that $F^k(\perp)$ is a **fixed point** of F .

$$F \begin{pmatrix} LVOut_1 \\ LVOut_2 \\ LVOut_3 \\ LVOut_4 \\ LVIn_1 \\ LVIn_2 \\ LVIn_3 \\ LVIn_4 \end{pmatrix} = \begin{pmatrix} LVIn_2 \\ LVIn_3 \cup LVIn_4 \\ LVIn_2 \\ \emptyset \\ LVOut_1 - \{x\} \\ LVOut_2 \cup \{y\} \\ LVOut_3 \cup \{x\} \\ LVOut_4 - \{x\} \end{pmatrix}$$

We get the following chain:

$$\underbrace{\begin{pmatrix} \emptyset \\ \emptyset \\ \emptyset \\ \emptyset \\ \emptyset \\ \emptyset \\ \emptyset \\ \emptyset \end{pmatrix}}_{\perp} \sqsubseteq \underbrace{\begin{pmatrix} \emptyset \\ \emptyset \\ \emptyset \\ \emptyset \\ \{y\} \\ \{x\} \\ \emptyset \end{pmatrix}}_{F(\perp)} \sqsubseteq \underbrace{\begin{pmatrix} \{y\} \\ \{x\} \\ \{y\} \\ \emptyset \\ \emptyset \\ \{y\} \\ \{x\} \\ \emptyset \end{pmatrix}}_{F^2(\perp)} \sqsubseteq \underbrace{\begin{pmatrix} \{y\} \\ \{x\} \\ \{y\} \\ \emptyset \\ \{y\} \\ \{x, y\} \\ \{x, y\} \\ \emptyset \end{pmatrix}}_{F^3(\perp)} \sqsubseteq \underbrace{\begin{pmatrix} \{x, y\} \\ \{x, y\} \\ \{x, y\} \\ \emptyset \\ \{y\} \\ \{x, y\} \\ \{x, y\} \\ \emptyset \end{pmatrix}}_{F^4(\perp)} = \underbrace{\begin{pmatrix} \{x, y\} \\ \{x, y\} \\ \{x, y\} \\ \emptyset \\ \{y\} \\ \{x, y\} \\ \{x, y\} \\ \emptyset \end{pmatrix}}_{F^5(\perp)}$$

Theorem

Assume a monotonically increasing function $F : L \rightarrow L$ and the following ascending Kleene chain:

$$\perp \sqsubseteq F(\perp) \sqsubseteq F^2(\perp) \sqsubseteq F^3(\perp) \sqsubseteq F^4(\perp) \sqsubseteq \dots$$

that stabilizes after the k -th iteration. Then:

- $F^k(\perp)$ is a fixed point of F .
- $F^k(\perp)$ is the least fixed point of F .
- Actually this result can be generalized to those cases in which L is infinite, and the chain does not stabilize.
- It requires stronger conditions on F (Scott-continuity).

SOLVING SYSTEMS OF EQUATIONS

- Assume that the variables in a system can take values from a **finite** poset.
1. Does F have a fixed point?
 - Yes, because the chain eventually stabilizes.
 - The chain will stabilize in a fixed point.
 2. Could it have several fixed points?
 - Yes.
 3. In case it has several fixed points, is there a fixed point that is better than the others?
 - The lower the fixed point, the more accurate it is, from the live variables analysis point-of-view.
 4. In that case, how can we compute the best fixed point?
 - By applying F successively on \perp until the chain stabilizes.

SOLVING SYSTEMS OF EQUATIONS

1. Does the system have a solution?
 - Yes, because the chain eventually stabilizes.
 - The chain will stabilize in a solution.
2. Could it have several solutions?
 - Yes.
3. In case it has several solutions, is there a solution that is better than the others?
 - The lower, the more accurate, from the live variables analysis point-of-view.
4. In that case, how can we compute the best solution?
 - By applying F successively on \perp until the chain stabilizes.

BACK TO OUR EXAMPLE

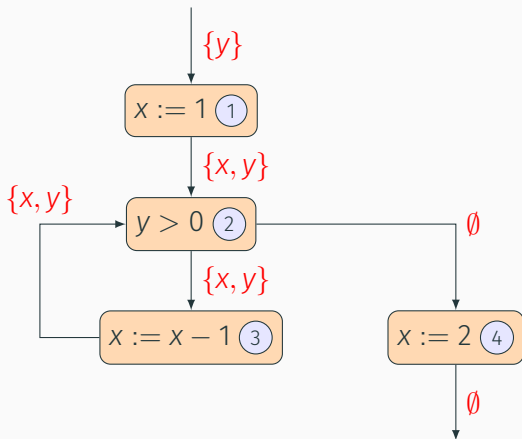
- In our example, the chain has stabilized at the following fixed point:

$$\underbrace{\begin{pmatrix} \{x, y\} \\ \{x, y\} \\ \{x, y\} \\ \emptyset \\ \{y\} \\ \{x, y\} \\ \{x, y\} \\ \emptyset \end{pmatrix}}_{F^5(\perp)}$$

which corresponds to the solution shown previously:

$$\begin{array}{llll} LVOut_1 = \{x, y\} & LVOut_2 = \{x, y\} & LVOut_3 = \{x, y\} & LVOut_4 = \emptyset \\ LVIn_1 = \{y\} & LVIn_2 = \{x, y\} & LVIn_3 = \{x, y\} & LVIn_4 = \emptyset \end{array}$$

EXAMPLE: OUR RESULT

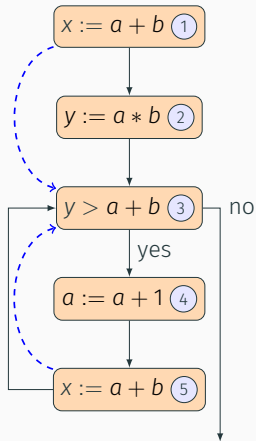


AVAILABLE EXPRESSIONS ANALYSIS

- We say that an expression e is **available** at a given program point p if all paths from the beginning of the program to p compute this expression, and the computation is not followed by an assignment to one of the operands in the expression.
- In other words, it determines those expressions that have been computed previously and need not be recomputed at p .
- An available expressions analysis determines, for each program point, which expressions **must** be available at that point.

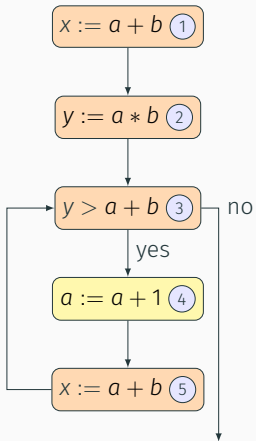
- **Common subexpression elimination** (CSE)
 - If an expression is computed in a given block, but it is already available at the entry of that block, it does not have to be recomputed.
 - Redundant expressions are stored in a temporary variable, and they are replaced by that variable in when they are available.

EXAMPLE



- Expression $a + b$ is available at block 3

EXAMPLE



- Expression $a + b$ is available at block 3
- Expression $a + b$ is not available at block 5, since the value of a is changed before reaching block 5.

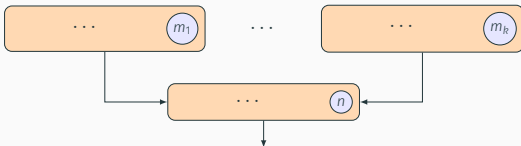
AVAILABLE EXPRESSIONS ANALYSIS

- We denote by **AExp** the set of nontrivial arithmetic expressions in the program.
 - Nontrivial = excluding variables and constants.
 - In our example: **AExp** = $\{a + b, a * b, a + 1\}$.
- The analysis computes, for each block, a pair of subsets of **AExp** which contain the available expressions at the entry and exit of the block, respectively.



AVAILABLE EXPRESSION ANALYSIS: DATA-FLOW EQUATIONS

- This time we will be doing a **forward analysis**.
 - An expression is available at the **entry** of a block if it is at the exit of **all** its predecessors.



$$AEIn_n = \bigcap_{m \in \text{pred}(n)} AEOut_m$$

- If n is the initial block, then $AEIn_n = \emptyset$.

- An expression is available at the **exit** of a block if:
 - It is computed in that block, and none of its operands is modified in it, or
 - It was available at the beginning of the block, and none of its operators is modified in it.

AVAILABLE EXPRESSION ANALYSIS: DATA-FLOW EQUATIONS

- For each block n , let us define $Kill_n$ as follows:
 - If the block contains an expression of the form $x := e$:

$$Kill_n = \{e' \in \mathbf{AExp} \mid x \text{ appears in } e'\}$$

- Otherwise, $Kill_n = \emptyset$.
- For each block n , let us define Gen_n as follows:
 - If n contains an expression of the form $x := e$:

$$Gen_n = \{e' \in \mathbf{AExp} \mid e' \text{ is a subexpression of } e \text{ not containing } x\}$$

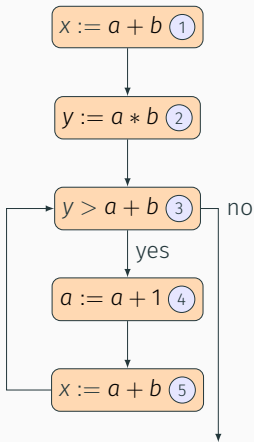
- If n contains a boolean expression:

$$Gen_n = \{e' \in \mathbf{AExp} \mid e' \text{ is a subexpression of } e\}$$

- Otherwise, $Gen_n = \emptyset$.

$$\text{Then: } AEOut_n = (AEIn_n - Kill_n) \cup Gen_n$$

AVAILABLE EXPRESSION ANALYSIS: DATA-FLOW EQUATIONS



$$Kill_1 = \emptyset$$

$$Kill_2 = \emptyset$$

$$Kill_3 = \emptyset$$

$$Kill_4 = \{a + b, a * b, a + 1\}$$

$$Kill_5 = \emptyset$$

$$Gen_1 = \{a + b\}$$

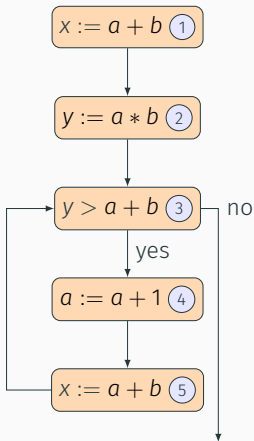
$$Gen_2 = \{a * b\}$$

$$Gen_3 = \{a + b\}$$

$$Gen_4 = \emptyset$$

$$Gen_5 = \{a + b\}$$

AVAILABLE EXPRESSION ANALYSIS: DATA-FLOW EQUATIONS



$$AEIn_1 = \emptyset$$

$$AEIn_2 = AEOut_1$$

$$AEIn_3 = AEOut_2 \cap AEOut_5$$

$$AEIn_4 = AEOut_3$$

$$AEIn_5 = AEOut_4$$

$$AEOut_1 = AEIn_1 \cup \{a + b\}$$

$$AEOut_2 = AEIn_2 \cup \{a * b\}$$

$$AEOut_3 = AEIn_3 \cup \{a + b\}$$

$$AEOut_4 = AEIn_4 - \{a + b, a * b, a + 1\}$$

$$AEOut_5 = AEIn_5 \cup \{a + b\}$$

- The unknowns of the equation system take values in $\mathcal{P}(\mathbf{AExp})$.
- A solution to the system is a vector \mathbf{AE} from $\mathcal{P}(\mathbf{AExp})^{10}$.
- Let us define an order relation \sqsubseteq in the same way as in the previous analysis.



ASSUME WE HAVE TWO SOLUTIONS \mathbf{AE} AND \mathbf{AE}' SUCH THAT $\mathbf{AE} \sqsubseteq \mathbf{AE}'$. WHICH ONE GIVES US MORE INFORMATION ON AVAILABLE EXPRESSIONS?

SOLVING THE SYSTEM OF EQUATIONS

- Let us define $F : \mathcal{P}(\text{AExp})^{10} \rightarrow \mathcal{P}(\text{AExp})^{10}$ as follows:

$$F \begin{pmatrix} \text{AEIn}_1 \\ \text{AEIn}_2 \\ \text{AEIn}_3 \\ \text{AEIn}_4 \\ \text{AEIn}_5 \\ \text{AEOut}_1 \\ \text{AEOut}_2 \\ \text{AEOut}_3 \\ \text{AEOut}_4 \\ \text{AEOut}_5 \end{pmatrix} = \begin{pmatrix} \emptyset \\ \text{AEOut}_1 \\ \text{AEOut}_2 \cap \text{AEOut}_5 \\ \text{AEOut}_3 \\ \text{AEOut}_4 \\ \text{AEIn}_1 \cup \{a + b\} \\ \text{AEIn}_2 \cup \{a * b\} \\ \text{AEIn}_3 \cup \{a + b\} \\ \text{AEIn}_4 - \{a + b, a * b, a + 1\} \\ \text{AEIn}_5 \cup \{a + b\} \end{pmatrix}$$

- This is an monotonically increasing function.
- The difference with the previous analysis is that now we are looking for a **greatest fixed point**.

DESCENDING KLEENE CHAIN

Theorem

Assume a monotonically increasing $F : L \rightarrow L$ and a Kleene descending chain

$$\top \supseteq F(\top) \supseteq F^2(\top) \supseteq F^3(\top) \supseteq F^4(\top) \supseteq \dots$$

that stabilizes at k -th iteration. Then:

- $F^k(\top)$ is a fixed point of F .
- $F^k(\top)$ is greater than any other fixed point of F .
- Therefore, in our available expression analysis, we have to apply successively F starting from \top . In this case:

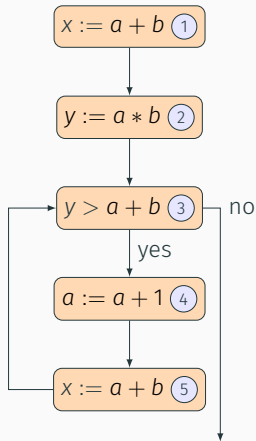
$$\top = \underbrace{(\text{AExp}, \text{AExp}, \dots, \text{AExp})}_{10 \text{ times}}$$

$$\begin{array}{c}
\left(\begin{array}{c} \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \end{array} \right) \rightarrow \left(\begin{array}{c} \emptyset \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \emptyset \\ \{a+b, a*b, a+1\} \end{array} \right) \rightarrow \left(\begin{array}{c} \emptyset \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \emptyset \\ \{a+b\} \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \emptyset \\ \{a+b, a*b, a+1\} \end{array} \right) \rightarrow \dots
\end{array}$$

$$\begin{array}{c}
\dots \rightarrow \left(\begin{array}{c} \emptyset \\ \{a+b\} \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \emptyset \\ \{a+b\} \\ \{a+b, a*b, a+1\} \\ \{a+b, a*b, a+1\} \\ \emptyset \\ \{a+b\} \end{array} \right) \rightarrow \left(\begin{array}{c} \emptyset \\ \{a+b\} \\ \{a+b\} \\ \{a+b, a*b, a+1\} \\ \emptyset \\ \{a+b\} \\ \{a+b, a*b\} \\ \{a+b, a*b, a+1\} \\ \emptyset \\ \{a+b\} \end{array} \right) \rightarrow \left(\begin{array}{c} \emptyset \\ \{a+b\} \\ \{a+b\} \\ \{a+b, a*b, a+1\} \\ \emptyset \\ \{a+b\} \\ \{a+b, a*b\} \\ \{a+b\} \\ \emptyset \\ \{a+b\} \end{array} \right) \rightarrow \dots
\end{array}$$

$$\dots \rightarrow \left(\begin{array}{c} \emptyset \\ \{a+b\} \\ \{a+b\} \\ \{a+b\} \\ \emptyset \\ \{a+b\} \\ \{a+b, a*b\} \\ \{a+b\} \\ \emptyset \\ \{a+b\} \end{array} \right) = \left(\begin{array}{c} \emptyset \\ \{a+b\} \\ \{a+b\} \\ \{a+b\} \\ \emptyset \\ \{a+b\} \\ \{a+b, a*b\} \\ \{a+b\} \\ \emptyset \\ \{a+b\} \end{array} \right)$$

ANALYSIS RESULTS



$$AEIn_1 = \emptyset$$

$$AEOut_1 = \{a + b\}$$

$$AEIn_2 = \{a + b\}$$

$$AEOut_2 = \{a + b, a * b\}$$

$$AEIn_3 = \{a + b\}$$

$$AEOut_3 = \{a + b\}$$

$$AEIn_4 = \{a + b\}$$

$$AEOut_4 = \emptyset$$

$$AEIn_5 = \emptyset$$

$$AEOut_5 = \{a + b\}$$

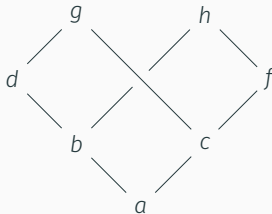
MONOTONE FRAMEWORKS

- To find a **general framework** able to express a vast amount of data-flow analyses.
- The data-flow analyses seen so far compute some values In_n and Out_n for each program block.
- But the kind of values computed depends on the particular analysis:
 - Live variables: $\mathcal{P}(\text{Var})$.
 - Available expressions: $\mathcal{P}(\text{AExp})$
- Although both compute values in $\mathcal{P}(X)$ for some X , this is not always the case.
 - For example, a constant propagation analysis returns a function $\text{Var} \rightarrow \mathbb{Z} \cup \{\perp, \top\}$ in each program point.

- In general, let us assume that In_n and Out_n take values in a set L .
- L is the **property space** of our analysis.
- In order to build an ascending or descending chain, L is required to have an order relation \sqsubseteq .
- L has to be a **lattice** so that the data-flow equations make sense.

LATTICES: UPPER BOUNDS AND LUBS

- Assume an ordered set (L, \sqsubseteq) :

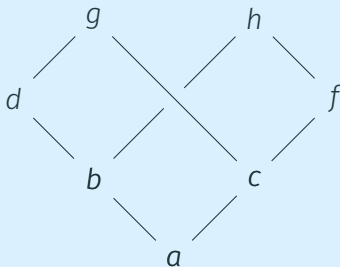


- Given a subset $L_0 \subseteq L$, an element $x \in L$ is an **upper bound** of L_0 if it is greater or equal than all the elements of L_0 :

$$\forall y \in L_0 : y \sqsubseteq x$$

- The **least upper bound** (lub) of L_0 , denoted by $\bigsqcup L_0$ is the lowest upper bound of all (if there is such lowest).

ASSUME THE FOLLOWING ORDERED SET L



• Upper bounds of $\{a, b, c\}$: $\sqcup\{a, b, c\}$

• Upper bounds of $\{a, c\}$: $\sqcup\{a, c\}$

• Upper bounds of $\{d, f\}$: $\sqcup\{d, f\}$

- Similarly, an element $x \in L$ is a **lower bound** of L_0 if it is lower or equal than all the elements of L_0 .

$$\forall y \in L_0 : x \sqsubseteq y$$

- The **greatest lower bound** (glb) of $L_0 \subseteq L$, denoted by $\sqcap L_0$ is the greatest lower bound of all (if there is such).
- When L_0 has two elements, we use $x \sqcup y$ instead of $\sqcup\{x, y\}$ and $x \sqcap y$ instead of $\sqcap\{x, y\}$.
- \sqcup and \sqcap are commutative, associative and idempotent.

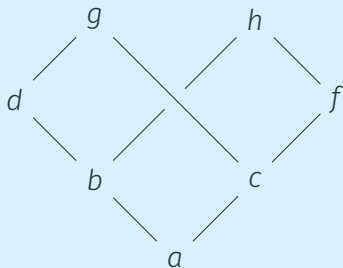
- A **lattice** is an ordered set (L, \subseteq) such that every pair of elements has a glb and lub in L .

$$\forall x, y \in L: x \sqcup y \in L \quad y \sqcap x \in L$$

- A **complete lattice** is an ordered set (L, \subseteq) in which every subset $L_0 \subseteq L$ has a glb and lub in L .

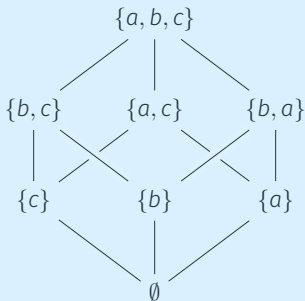
$$\forall L_0 \subseteq L. \bigsqcup L_0 \in L \quad y \quad \bigsqcap L_0 \in L$$

IS THE FOLLOWING ORDERED SET A LATTICE?



.

Is $(\mathcal{P}(\{a, b, c\}), \subseteq)$ A LATTICE? AND A COMPLETE LATTICE?



.

- In general, for every set X , the ordered set $(\mathcal{P}(X), \subseteq)$ is a complete lattice.
- In this case, \sqcup is set union (\cup) and \sqcap is set intersection (\cap).

Is (\mathbb{N}, \leq) A LATTICE? IS IT A COMPLETE LATTICE?



⋮
2
|
1
|
0

•

•

WHAT IF WE ADD A NEW ELEMENT $+\infty$?



- Assume that $x \leq +\infty$ para todo $x \in \mathbb{N}^\infty$.

$+\infty$
⋮
2
|
1
|
0

•

- Complete lattices are bounded sets. That is, they have \perp and \top elements.
- In fact:

$$\perp = \bigcap L \qquad \top = \bigcup L$$

- Given (L_1, \sqsubseteq_1) , (L_2, \sqsubseteq_2) , and an order relation \sqsubseteq on $L_1 \times L_2$.

Theorem

If L_1 and L_2 are (complete) lattices, so is $L_1 \times L_2$.

- Let \sqcup_1 and \sqcup_2 be the lub operators in (L_1, \sqsubseteq_1) and (L_2, \sqsubseteq_2) respectively. Then, for every $(x_1, x_2), (y_1, y_2) \in L_1 \times L_2$:

$$(x_1, x_2) \sqcup (y_1, y_2) = (x_1 \sqcup_1 y_1, x_2 \sqcup_2 y_2)$$

- Similarly with greatest lower bounds:

$$(x_1, x_2) \sqcap (y_1, y_2) = (x_1 \sqcap_1 y_1, x_2 \sqcap_2 y_2)$$

- This can be extended to n -ary products: $L_1 \times L_2 \times \cdots \times L_n$:

$$(x_1, x_2, \dots, x_n) \sqcup (y_1, y_2, \dots, y_n) = (x_1 \sqcup_1 y_1, x_2 \sqcup_2 y_2, \dots, x_n \sqcup_n y_n)$$

$$(x_1, x_2, \dots, x_n) \sqcap (y_1, y_2, \dots, y_n) = (x_1 \sqcap_1 y_1, x_2 \sqcap_2 y_2, \dots, x_n \sqcap_n y_n)$$

for every $(x_1, \dots, x_n), (y_1, \dots, y_n) \in L_1 \times \cdots \times L_n$.

- The \perp element in $L_1 \times \cdots \times L_n$ is $(\perp_1, \dots, \perp_n)$.
- The \top element in $L_1 \times \cdots \times L_n$ is (\top_1, \dots, \top_n) .

LATTICE CONSTRUCTION: TOTAL FUNCTIONS.

- Given a set S and an ordered set (L_1, \sqsubseteq_1) , we denote by $S \rightarrow L_1$ the set of total functions from S to L_1 .
- Let us define a relation \sqsubseteq in $S \rightarrow L_1$ as follows:
 - Given $f_1, f_2 : S \rightarrow L_1$, we say that $f_1 \sqsubseteq f_2$ if and only if:

$$\forall s \in S : f_1(s) \sqsubseteq_1 f_2(s)$$

Theorem

- \sqsubseteq is an order relation.
- If L_1 is a (complete) lattice, so is $(S \rightarrow L_1, \sqsubseteq)$.
- If L_1 is a lattice we have, for any $s \in S$:

$$(f_1 \sqcup f_2)(s) = f_1(s) \sqcup_1 f_2(s)$$

$$(f_1 \sqcap f_2)(s) = f_1(s) \sqcap_1 f_2(s)$$



WHICH ARE THE \top AND \perp ELEMENTS OF $S \rightarrow L_1$?

-
-

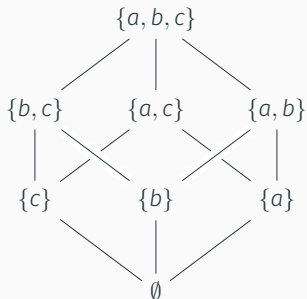
- Given an ordered set (L, \sqsubseteq) , let us consider the inverse relation \sqsupseteq .

Theorem

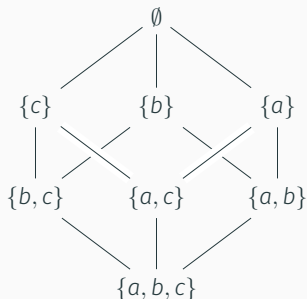
- \sqsupseteq is an order relation.
- If (L, \sqsubseteq) is a complete lattice, so is (L, \sqsupseteq) .
- If (L, \sqsubseteq) is a lattice:
 - The \sqcup operator in (L, \sqsupseteq) is the \sqcap operator in (L, \sqsubseteq) .
 - The \sqcap operator in (L, \sqsupseteq) is the \sqcup operator in (L, \sqsubseteq) .
 - The \perp element in (L, \sqsupseteq) is the \top element in (L, \sqsubseteq) .
 - The \top element in (L, \sqsupseteq) is the \perp element in (L, \sqsubseteq) .

LATTICE CONSTRUCTION: DUALITY

$(\mathcal{P}(\{a, b, c\}), \subseteq)$



$(\mathcal{P}(\{a, b, c\}), \supseteq)$



- Given a poset (L, \sqsubseteq) , an **ascending chain** (l_n) is a sequence $l_0, l_1, \dots, l_n, \dots$ of elements in L such that:

$$l_0 \sqsubseteq l_1 \sqsubseteq \dots \sqsubseteq l_n \sqsubseteq \dots$$

- We say that an ascending chain (l_n) **stabilizes** if all its elements are equal from a given one.

$$l_0 \sqsubset l_1 \sqsubset \dots \sqsubset l_k = l_{k+1} = \dots$$

- We say that a poset (L, \sqsubseteq) satisfies the **ascending chain condition** if every possible chain in L stabilizes.

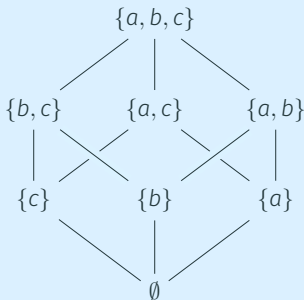
DOES $(\mathbb{N}^\infty, \leq)$ SATISFY THE ASCENDING
CHAIN CONDITION?



$+\infty$
⋮
2
|
1
|
0

•

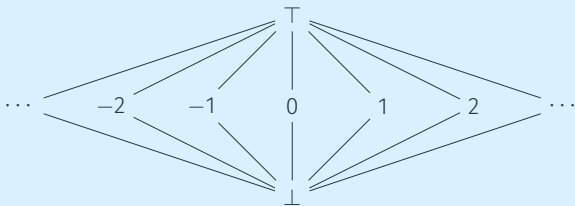
DOES $(\mathcal{P}(\{a, b, c\}), \subseteq)$ SATISFY THE
ASCENDING CHAIN CONDITION?



.



DOES THE FOLLOWING LATTICE
 $(\mathbb{Z} \cup \{\perp, \top\}, \sqsubseteq)$ SATISFY THE ASCENDING
CHAIN CONDITION?



.

Theorem

If L_1, L_2, \dots, L_n satisfy the ascending chain condition, so does $L_1 \times \dots \times L_n$.

Theorem

If L satisfies the ascending chain condition and S is finite, $S \rightarrow L$ also satisfies the ascending chain condition.

- A **monotone framework** is made up of:
 - A complete lattice L satisfying the ascending chain condition.
 - We denote by \sqcup the lub operator associated with L .
 - A set \mathcal{F} of monotonically increasing functions $L \rightarrow L$ (**transfer functions**) such that:
 - \mathcal{F} contains the identity function.
 - \mathcal{F} is closed under function composition:

$$\forall f, g \in \mathcal{F} : f \circ g \in \mathcal{F}$$

- Assume we want to define a control-flow analysis in a given program.
- An **instance** of a monotone framework is made up of:
 - A monotone framework (L, \mathcal{F}) .
 - An **extremal value**, applied to the initial block of the program.
 - For every block n in the program, a **transfer function** $f_n \in \mathcal{F}$.

- From an instance in a monotone framework, we can define the equations of a forward analysis as follows:

$$In_n = \begin{cases} \iota & \text{if } n \text{ is the initial block} \\ \sqcup \{Out_m \mid m \in pred(n)\} & \text{otherwise} \end{cases}$$

$$Out_n = f_n(In_n)$$

BACKWARD ANALYSIS: GENERAL DEFINITION

- The equations of a backward analysis are defined as follows:

$$Out_n = \begin{cases} \iota & \text{if } n \text{ is a final block} \\ \sqcup \{In_m \mid m \in succ(n)\} & \text{otherwise} \end{cases}$$

$$In_n = f_n(Out_n)$$

- Some formalizations unify the equations of forward and backward analyses.
 - In order to perform a backward analysis, it is enough to reverse the direction of the arrows of the flow graph, and perform a forward analysis in the resulting program.

MONOTONE FRAMEWORKS: OUR ANALYSES

- Lattices used so far:

	<i>Live. variables</i>	<i>Reaching defs.</i>	<i>Available exps.</i>
L	$\mathcal{P}(\text{Var})$	$\mathcal{P}(\text{Var} \times \text{Lab}^?)$	$\mathcal{P}(\text{AExp})$
\sqsubseteq	\subseteq	\subseteq	\supseteq
\sqcup	\cup	\cup	\cap
\perp	\emptyset	\emptyset	AExp
ι	\emptyset	$\{(x, ?) \mid x \in \text{Var}\}$	\emptyset

- In these cases, the set \mathcal{F} of transfer functions is made up of those functions $f: L \rightarrow L$ such that for every $l \in L$ there exist l_k and l_g such that:

$$f(l) = (l - l_k) \cup l_g$$



F. Nielson, H.R. Nielson, C. Hankin

Principles of Program Analysis

Springer (1999)

Chapter 2



U.P. Khedker, A. Sanyal, B. Karkare

Data Flow Analysis: Theory and Practice

CRC Press (2009)

