

Principles of Abstract Interpretation

MIT press

Ch. 41, Dataflow Analysis

Patrick Cousot

pcousot.github.io

PrAbsInt@gmail.com github.com/PrAbsInt/

These slides are available at
<http://github.com/PrAbsInt/slides/slides-41--data-flow-analysis-PrAbsInt.pdf>

Ch. 41, Dataflow Analysis

(Classic) data flow analysis

- **Dataflow analysis**: at each point of a computer program gather information about how variables and/or expressions are used and modified
- Method:
 - Sets (of variables, expressions, etc.) are finite and represented by bit vectors
 - Program \rightarrow control flow graph \rightarrow equations \rightarrow fixpoint computation \rightarrow solution
- Mainly used for compilation
- Examples:
 - **Reaching definitions**: which assignments to x reach a use of a variable y in an expression at a program point without another assignments to x in between
 - **Available expressions**: which expressions previously computed expressions need not be recomputed at a program point (since its variables are unchanged)
 - **Live variable analysis**: which variables hold a value that may be used later
 - ...

Soundness of a data flow analysis

- Originally **postulated equationally** [Allen & Cocke, NYU/IBM, 70's]
 - **Postulated on paths** in the control flow graph to justify the equations [Kildall, 1973]
 - The dataflow analysis on a path can be **specified by a temporal logic formula** (considering the control flow graph as a transition system) [Steffen, 1991, 1993]
 - The **control flow chart** is an abstract interpretation of the program semantics [Schmidt, 1998]
- The soundness is not with respect to the program semantics but with respect to an abstraction of the program semantics (i.e. the control flow chart)
- This is **problematic**
- We study a **correct alternative**: soundness by abstract interpretation with respect to a semantics [P. Cousot and R. Cousot, 1979, Section 7.2.0.6.3 *Justifying the Data Flow Equations of "Available Expressions"*]

Live and dead variables analysis

- A variable x is **live** at a program point on a program path iff it may be used before being modified.
- A variable x is **dead** at a program point on a program path iff it will definitely be modified before being used

ℓ_1	\longleftarrow	x and y dead (modified in ℓ_1 and ℓ_2 before being used)
$x = 1 ;$		
ℓ_2	\longleftarrow	x live (used in expression at ℓ_2) and y dead (modified in ℓ_2 before use)
$y = x ;$		
ℓ_3	\longleftarrow	x dead (modified in ℓ_3 before use) and y live (used in expression at ℓ_3)
$x = y + 1 ;$		
ℓ_4	\longleftarrow	x and y live (used in assigned expression at ℓ_4)
$x = x - y ;$		
ℓ_5	\longleftarrow	x live and y dead (by hypothesis)
$L_e = \{x\}$	\longleftarrow	hypothesis that x is live and y is dead on statement exit.

Liveness, informally

“a variable is **potentially**/**definitely** live at some program point ℓ if it holds a value that **may**/**must** be used in the future before the next time the variable is modified”.

The liveness abstraction of a trace

↓ This l is for “liveness”

$\alpha_{use,mod}^l \llbracket S \rrbracket L_b, L_e \langle \pi_0^\ell, {}^\ell\pi \rangle$

- ${}^\ell$ is the program label ${}^\ell = \text{at} \llbracket S \rrbracket$
- π_0^ℓ is an initial trace of the program component S arriving ${}^\ell = \text{at} \llbracket S \rrbracket$
- ${}^\ell\pi$ continues the initial trace π_0 from ${}^\ell = \text{at} \llbracket S \rrbracket$ on
- use defines the set $use \llbracket a \rrbracket \rho$ of variables which value is used when executing action a in environment ρ
- mod defines the set $mod \llbracket a \rrbracket \rho$ of variables which value is modified when executing action a in environment ρ .
- L_b is the set of variables live on exit of S by a **break** ;, if any
- L_e is the set of variables live on exit of S by a normal exit, if ever

$\alpha_{use,mod}^l \llbracket S \rrbracket L_b, L_e \langle \pi_0^\ell, {}^\ell\pi \rangle$ is the set of variables live at ${}^\ell$ in ${}^\ell\pi$ continuing π_0^ℓ .

The liveness abstraction of a trace, formally & recursively

$$\alpha_{use,mod}^l \llbracket S \rrbracket L_b, L_e \langle \pi_0, \ell \rangle \triangleq \{x \in \mathbb{V} \mid (\ell = \text{after} \llbracket S \rrbracket \wedge x \in L_e) \vee \quad (a_1) \quad (41.2)$$

$$(\text{escape} \llbracket S \rrbracket \wedge \ell = \text{break-to} \llbracket S \rrbracket \wedge x \in L_b)\} \quad (a_2)$$

$$\alpha_{use,mod}^l \llbracket S \rrbracket L_b, L_e \langle \pi_0, \ell \xrightarrow{a} \ell' \pi_1 \rangle \triangleq \{x \in \mathbb{V} \mid x \in use \llbracket a \rrbracket \varrho(\pi_0) \vee \quad (b_1)$$

$$(x \notin mod \llbracket a \rrbracket \varrho(\pi_0) \wedge x \in \alpha_{use,mod}^l \llbracket S \rrbracket L_b, L_e \langle \pi_0 \frown \ell \xrightarrow{a} \ell', \ell' \pi_1 \rangle)\} \quad (b_2)$$

The liveness abstraction of a trace, formally & recursively

$$\alpha_{use,mod}^l \llbracket S \rrbracket L_b, L_e \langle \pi_0, \ell \rangle \triangleq \{x \in \mathcal{V} \mid (\ell = \text{after} \llbracket S \rrbracket \wedge x \in L_e) \vee \quad (a_1) \quad (41.2)$$

$$(\text{escape} \llbracket S \rrbracket \wedge \ell = \text{break-to} \llbracket S \rrbracket \wedge x \in L_b)\} \quad (a_2)$$

$$\alpha_{use,mod}^l \llbracket S \rrbracket L_b, L_e \langle \pi_0, \ell \xrightarrow{a} \ell' \pi_1 \rangle \triangleq \{x \in \mathcal{V} \mid x \in use \llbracket a \rrbracket \varrho(\pi_0) \vee \quad (b_1)$$

$$(x \notin mod \llbracket a \rrbracket \varrho(\pi_0) \wedge x \in \alpha_{use,mod}^l \llbracket S \rrbracket L_b, L_e \langle \pi_0, \ell \xrightarrow{a} \ell', \ell' \pi_1 \rangle)\} \quad (b_2)$$

ℓ_1	\longleftarrow	$\neg(b_1) \wedge \neg(b_2)$ x and y dead (modified in ℓ_1 and ℓ_2 before being used)
		x = 1 ;
ℓ_2	\longleftarrow	(b ₂) x live (used in expression at ℓ_3) and $\neg(b_1) \wedge \neg(b_2)$ y dead (modified in ℓ_2 before use)
		y = 1 ;
ℓ_3	\longleftarrow	(b ₁) x live (used in ℓ_3) and (b ₂) y live (not modified and used in expression at ℓ_4)
		x = x + 1 ;
ℓ_4	\longleftarrow	(b ₁) x and (b ₁) y live (used in assigned expression at ℓ_4)
		x = x - y ;
ℓ_5	\longleftarrow	(a ₁) x live and $\neg(a_1) \wedge \neg(a_2)$ y dead (by hypothesis)
		$L_e = \{x\}$ \longleftarrow hypothesis that x is live and y is dead on statement exit.

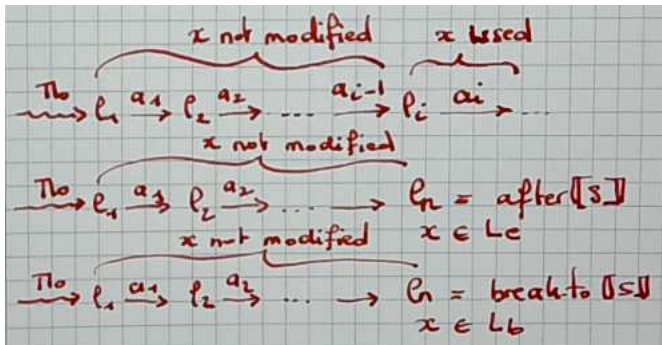
Potential and definite liveness of a statement, formally & iteratively

Lemma 1 If $\pi_1 = \ell_1 \xrightarrow{a_1} \ell_2 \xrightarrow{a_2} \dots \xrightarrow{a_{n-1}} \ell_n$ and $\langle \pi_0, \pi_1 \rangle \in \mathcal{S}^* \llbracket S \rrbracket$ then

$$\alpha_{use, mod}^l \llbracket S \rrbracket L_b, L_e \langle \pi_0, \pi_1 \rangle = \{x \in \mathcal{V} \mid \exists i \in [1, n-1] . \forall j \in [1, i-1] .$$

$$x \notin mod \llbracket a_j \rrbracket \varrho(\pi_0 \frown \ell_1 \xrightarrow{a_1} \ell_2 \dots \xrightarrow{a_{j-1}} \ell_j) \wedge x \in use \llbracket a_i \rrbracket \varrho(\pi_0 \frown \ell_1 \xrightarrow{a_1} \ell_2 \dots \xrightarrow{a_{i-1}} \ell_i)\}$$

$$\cup (\ell_n = after \llbracket S \rrbracket \text{ ? } L_e : \emptyset) \cup (escape \llbracket S \rrbracket \wedge \ell_n = break-to \llbracket S \rrbracket \text{ ? } L_b : \emptyset). \quad \square$$



Potential and definite liveness of a statement, formally

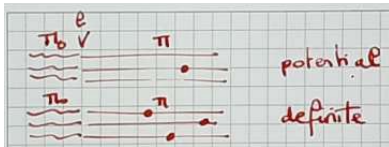
The potential and definite liveness are abstractions of the maximal trace semantics $\mathcal{S}^{+\infty}[[S]]$ is by merge over all traces

$$\alpha_{use,mod}^{\exists l}[[S]] \mathcal{S} L_b, L_e \triangleq \bigcup_{\langle \pi_0, \pi \rangle \in \mathcal{S}} \alpha_{use,mod}^l[[S]] L_b, L_e \langle \pi_0, \pi \rangle \quad \text{potential liveness} \quad (41.3)$$

$$\alpha_{use,mod}^{\forall l}[[S]] \mathcal{S} L_b, L_e \triangleq \bigcap_{\langle \pi_0, \pi \rangle \in \mathcal{S}} \alpha_{use,mod}^l[[S]] L_b, L_e \langle \pi_0, \pi \rangle \quad \text{definite liveness} \quad (41.4)$$

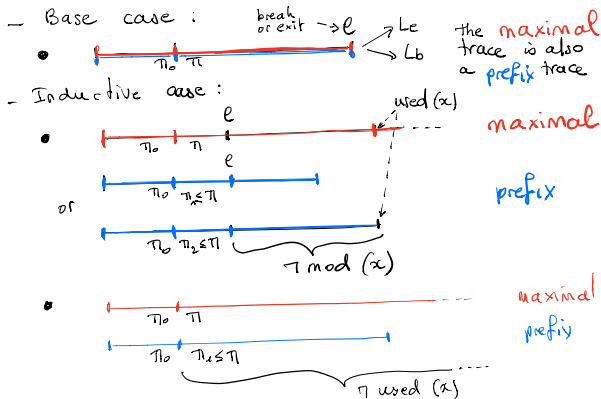
$$\alpha_{use,mod}^{\exists d}[[S]] \mathcal{S} D_b, D_e \triangleq \neg \alpha_{use,mod}^{\forall l}[[S]] \mathcal{S} \neg D_b, \neg D_e \quad \text{potential deadness} \quad (41.5)$$

$$\alpha_{use,mod}^{\forall d}[[S]] \mathcal{S} D_b, D_e \triangleq \neg \alpha_{use,mod}^{\exists l}[[S]] \mathcal{S} \neg D_b, \neg D_e \quad \text{definite deadness} \quad (41.6)$$



Prefix versus maximal trace semantics

Lemma 2 41.9 Using the prefix trace semantics $\mathcal{S}^* \llbracket S \rrbracket$ or the maximal trace semantics $\mathcal{S}^{+\infty} \llbracket S \rrbracket$ in the definition of potential liveness $\alpha_{use, mod}^{\exists l} \llbracket S \rrbracket$ is equivalent. □



Semantic liveness/deadness abstractions $\mathcal{S}^{\exists!} \llbracket S \rrbracket \triangleq \alpha_{\text{use}, \text{mod}}^{\exists!} \llbracket S \rrbracket$

- An action a uses variable y in a given environment ρ if and only if it is possible to change the value of y so as to change the effect of action a on program execution.

$$\text{use} \llbracket \text{skip} \rrbracket \rho \triangleq \emptyset \quad (41.10)$$

$$\text{use} \llbracket x = A \rrbracket \rho \triangleq \{y \mid \exists v \in \mathbb{V} . \mathcal{A} \llbracket A \rrbracket \rho \neq \mathcal{A} \llbracket A \rrbracket \rho[y \leftarrow v] \wedge \rho(x) \neq \mathcal{A} \llbracket A \rrbracket \rho\}$$

$$\text{use} \llbracket a \rrbracket \rho \triangleq \{y \mid \exists v \in \mathbb{V} . \mathcal{B} \llbracket a \rrbracket \rho \neq \mathcal{B} \llbracket a \rrbracket \rho[y \leftarrow v]\} \quad a \in \{B, \neg(B)\}$$

Example, $y \notin \text{use} \llbracket x = y - y \rrbracket \rho$ and $x \notin \text{use} \llbracket x = x \rrbracket \rho$.

- An action a modifies variable x in an environment ρ if and only if the execution of action a in environment ρ changes the value of x .

$$\text{mod} \llbracket a \rrbracket \rho \triangleq \{x \mid a = (x = A) \wedge (\rho(x) \neq \mathcal{A} \llbracket A \rrbracket \rho)\}$$

Classic syntactic liveness/deadness abstractions $\mathcal{S}^{\exists!} \llbracket S \rrbracket \triangleq \alpha_{\text{use}, \text{mod}}^{\exists!} \llbracket S \rrbracket$

The set $\text{use} \llbracket a \rrbracket$ of variables used and the set $\text{mod} \llbracket a \rrbracket$ of variables assigned to/modified in an action $a \in \mathcal{A}$ are postulated to be as follows (the parameter ρ is useless but added for consistency with (41.2)).

$$\text{use} \llbracket x = A \rrbracket \rho \triangleq \text{vars} \llbracket A \rrbracket$$

$$\text{mod} \llbracket x = A \rrbracket \rho \triangleq \{x\}$$

$$\text{use} \llbracket \text{skip} \rrbracket \rho \triangleq \emptyset$$

$$\text{mod} \llbracket \text{skip} \rrbracket \rho \triangleq \emptyset$$

$$\text{use} \llbracket B \rrbracket \rho \triangleq \text{use} \llbracket \neg(B) \rrbracket \rho \triangleq \text{vars} \llbracket B \rrbracket$$

$$\text{mod} \llbracket B \rrbracket \rho \triangleq \text{mod} \llbracket \neg(B) \rrbracket \rho \triangleq \emptyset$$

where $\text{vars} \llbracket E \rrbracket$ is the set of program variables occurring in arithmetic or boolean expression E .

Unsoundness of the syntactic liveness/deadness abstractions


$$\mathcal{S}^{\exists!}[[S]] \not\subseteq \mathcal{S}^{\exists!!}[[S]]$$

- Counter-example: $\ell_1 \ x = y - y ; \ell_2$ where x is live at ℓ_2 on exit
- Syntactically,
 - x is not used in $y-y$
 - x is modified by the assignment so x is always syntactically dead at ℓ_1 .
- Semantically,
 - x is not used in $y-y$ (since changing the value of x at ℓ_1 will not change the value of $y-y$ which is always 0).
 - x is not always modified by the assignment $x = y - y$; (in case x was 0 before)
- The problem is that

$$\begin{array}{l} \exists \rho \in \mathbb{E}_v . y \in \text{use}[[a]] \rho \Rightarrow \forall \rho \in \mathbb{E}_v . y \in \text{use}[[a]] \rho \\ \text{but} \quad \exists \rho \in \mathbb{E}_v . x \in \text{mod}[[a]] \rho \wedge x \notin \text{mod}[[a]] \rho \end{array}$$

What could go wrong when optimizing programs?

Consider a compiler that successively performs

1. a (syntactic) liveness analysis ;
2. next, a code optimization by removal
 - (a) of assignments to variables that are dead after this assignment,
 - (b) of assignments to variables that do not change the value of this variable (using Kildall's constancy analysis [Kildall, 1973] or better symbolic constancy analysis [Haghighat and Polychronopoulos, 1996; Wegman and Zadeck, 1991]);
3. next, a register allocation such that
 - (a) simultaneously live variables are stored in different registers,
 - (b) when no register is left and one is needed, one of those containing the value of a dead variable is preferred (to avoid saving the value of the variable to its memory location as would be needed for live variables).

For the following program (where all variables are dead on exit)

	semantically		syntactically		
	live	dead	live	dead	
x=0; scanf(y);					
if (x==0){					
ℓ ₁ ... x and y neither used nor modified ...	ℓ ₁	{x}	{y}	{y}	{x}
ℓ ₂ x = y - y; }	ℓ ₂	{x}	{y}	{y}	{x}
else {					
x=42;					
}					
ℓ ₃ print(x);	ℓ ₃	{x}	{y}	{x}	{y}

- Code elimination ((2.b)) suppresses the assignment at ℓ_2 since the value of x is unchanged.
- Assume x is in a register at ℓ_1 and a fresh register is needed but none is left available. By ((3.b)) the register containing x may be selected since its value need not be saved to memory because x is syntactically dead at ℓ_1 .
- Then the value of x is lost at ℓ_3 , a compilation bug.
- This error does not occur with semantic liveness.

Solutions not to go wrong

- Prevent program transformations (such as (2.b) and (3.b) above) that do not preserve the soundness of the semantic liveness $\mathcal{S}^{\exists!}$.
- Move elimination of assignments to variables that do not change the value of this variable ((2.b)) before liveness analysis.
- Redo the liveness analysis after any program transformation that does not preserve the information.
- A better solution is adopted in CompCert [Leroy, 2009]: the liveness analysis and code elimination are performed simultaneously and the liveness analysis is updated to be valid *after* code elimination.

Computational design of the structural syntactic potential liveness analysis

$$\begin{aligned}
 \hat{\mathcal{S}}^{\exists}[\text{S}\ell] L_e &\triangleq \hat{\mathcal{S}}^{\exists}[\text{S}\ell] \emptyset, L_e & (41.22) \\
 \hat{\mathcal{S}}^{\exists}[\text{x} = \text{A} ;] L_b, L_e &\triangleq \text{use}[\text{x} = \text{A}] \cup (L_e \setminus \text{mod}[\text{x} = \text{A}]) \\
 \hat{\mathcal{S}}^{\exists}[;] L_b, L_e &\triangleq L_e \\
 \hat{\mathcal{S}}^{\exists}[\text{S}\ell' \text{ S}] L_b, L_e &\triangleq \hat{\mathcal{S}}^{\exists}[\text{S}\ell'] L_b, (\hat{\mathcal{S}}^{\exists}[\text{S}] L_b, L_e) \\
 \hat{\mathcal{S}}^{\exists}[\epsilon] L_b, L_e &\triangleq L_e \\
 \hat{\mathcal{S}}^{\exists}[\text{if}(\text{B}) \text{S}_t] L_b, L_e &\triangleq \text{use}[\text{B}] \cup L_e \cup \hat{\mathcal{S}}^{\exists}[\text{S}_t] L_b, L_e \\
 \hat{\mathcal{S}}^{\exists}[\text{if}(\text{B}) \text{S}_t \text{ else } \text{S}_f] L_b, L_e &\triangleq \text{use}[\text{B}] \cup \hat{\mathcal{S}}^{\exists}[\text{S}_t] L_b, L_e \cup \hat{\mathcal{S}}^{\exists}[\text{S}_f] L_b, L_e \\
 \hat{\mathcal{S}}^{\exists}[\text{while}(\text{B}) \text{S}_b] L_b, L_e &\triangleq \text{use}[\text{B}] \cup L_e \cup \hat{\mathcal{S}}^{\exists}[\text{S}_b] L_b, L_e \\
 \hat{\mathcal{S}}^{\exists}[\text{break} ;] L_b, L_e &\triangleq L_b \\
 \hat{\mathcal{S}}^{\exists}[\{\text{S}\ell\}] L_b, L_e &\triangleq \hat{\mathcal{S}}^{\exists}[\text{S}\ell] L_b, L_e
 \end{aligned}$$

No fixpoint iteration for the `while`, the solution can be directly computed by 1 iteration!.

Theorem 41.24 $\hat{\mathcal{S}}^{\exists l} \llbracket S \rrbracket$ defined by (41.22) is syntactically sound that is

$$\alpha_{\text{use}, \text{mod}}^{\exists l} \not\subseteq \alpha_{\text{use}, \text{mod}}^{\exists l} (\mathcal{S}^* \llbracket S \rrbracket) \subseteq \alpha_{\text{use}, \text{mod}}^{\exists l} \llbracket S \rrbracket (\mathcal{S}^* \llbracket S \rrbracket) \subseteq \hat{\mathcal{S}}^{\exists l} \llbracket S \rrbracket$$

Informally:

- A variable is live at a program component iff it is not semantically/syntactically used before being (syntactically) *assigned to*
- So, a compiler removing this assignment would be incorrect.

Structural syntactic definite deadness analysis

$$\hat{\mathcal{S}}^{\forall d}[[S]] D_b, D_e \triangleq \neg \hat{\mathcal{S}}^{\exists d}[[S]] \neg D_b, \neg D_e$$

$$\begin{aligned} \hat{\mathcal{S}}^{\forall d}[[S \downarrow \ell]] D_e &= \hat{\mathcal{S}}^{\forall d}[[S \downarrow \ell]] \vee, D_e & (41.6) \\ \hat{\mathcal{S}}^{\forall d}[[x = A;]] D_b, D_e &= \neg \text{use}[[x = A]] \cap (D_e \cup \text{mod}[[x = A]]) \\ \hat{\mathcal{S}}^{\forall d}[[;]] D_b, D_e &= D_e \\ \hat{\mathcal{S}}^{\forall d}[[S \downarrow' S]] D_b, D_e &= \hat{\mathcal{S}}^{\forall d}[[S \downarrow']] D_b, (\hat{\mathcal{S}}^{\forall d}[[S]] D_b, D_e) \\ \hat{\mathcal{S}}^{\forall d}[[\epsilon]] D_b, D_e &= D_e \\ \hat{\mathcal{S}}^{\forall d}[[\text{if } (B) S_t]] D_b, D_e &= \neg \text{use}[[B]] \cap D_e \cap \hat{\mathcal{S}}^{\forall d}[[S_t]] D_b, D_e \\ \hat{\mathcal{S}}^{\forall d}[[\text{if } (B) S_t \text{ else } S_f]] D_b, D_e &= \neg \text{use}[[B]] \cap \hat{\mathcal{S}}^{\forall d}[[S_t]] D_b, D_e \cap \hat{\mathcal{S}}^{\forall d}[[S_f]] D_b, D_e \\ \hat{\mathcal{S}}^{\forall d}[[\text{while } (B) S_b]] D_b, D_e &= \neg \text{use}[[B]] \cap D_e \cap \hat{\mathcal{S}}^{\forall d}[[S_b]] D_b, D_e \\ \hat{\mathcal{S}}^{\forall d}[[\text{break};]] D_b, D_e &= D_b \\ \hat{\mathcal{S}}^{\forall d}[[\{ S \downarrow \}]] D_b, D_e &= \hat{\mathcal{S}}^{\forall d}[[S \downarrow]] D_b, D_e \end{aligned}$$

Conclusion

- Correct compilation should preserve the source semantics
- Using syntactic reasonings may be problematic if not well-understood with respect to the source, intermediate, and object code semantics
- The tradition of using a control graph, boolean equations, and fixpoint iteration may be inefficient
- Structural induction is more efficient (and no compiler writer looks to know that!)

References I

Bibliography

- Allen, Frances E. (1970). "Control Flow Analysis.". *SIGPLAN Not.* 5.7, pp. 1–19.
- (1971). "A Basis for Program Optimization.". In *IFIP Congress (1)*. Pp. 385–390.
- (1974). "Interprocedural data flow analysis.". In Jack L. Rosenfeld, ed. *Information Processing 74*. North-Holland Pub. Co., pp. 398–402.
- Allen, Frances E. and John Cocke (1976). "A Program Data Flow Analysis Procedure.". *Commun. ACM*. 19.3, pp. 137–147.
- Cousot, Patrick and Radhia Cousot (1979). "Systematic Design of Program Analysis Frameworks.". In *POPL*. ACM Press, pp. 269–282.
- Haghighat, Mohammad R. and Constantine D. Polychronopoulos (1996). "Symbolic Analysis for Parallelizing Compilers.". *ACM Trans. Program. Lang. Syst.* 18.4, pp. 477–518.

References II

- Kildall, Gary A. (1973). "A Unified Approach to Global Program Optimization." In *POPL*. ACM Press, pp. 194–206.
- Leroy, Xavier (2009). "Formal verification of a realistic compiler." *Commun. ACM*. 52.7, pp. 107–115.
- Schmidt, David A. (1998). "Data Flow Analysis is Model Checking of Abstract Interpretations." In *POPL*. ACM, pp. 38–48.
- Steffen, Bernhard (1991). "Data Flow Analysis as Model Checking." In *TACS*. Vol. 526. Lecture Notes in Computer Science. Springer, pp. 346–365.
- (1993). "Generating Data Flow Analysis Algorithms from Modal Specifications." *Sci. Comput. Program.* 21.2, pp. 115–139.
- Wegman, Mark N. and F. Kenneth Zadeck (1991). "Constant Propagation with Conditional Branches." *ACM Trans. Program. Lang. Syst.* 13.2, pp. 181–210.

Home work

Read Ch. 41 “Dataflow Analysis” of

Principles of Abstract Interpretation

Patrick Cousot

MIT Press

The End, Thank you