# Lecture 12:
# Data flow analysis

John Wickerson

Compilers

# Anatomy of a compiler

```
fu = a + a - zo
```

source code (stream of characters)

```
f u  =  a  +  a  -  z o
```

**lexer**

stream of tokens

```
fu = a + a - zo
```

**parser**

syntax tree

```
        assign
       /      \
     fu       sub
             /    \
           add     zo
          /   \
         a     a
```

**type checker**

syntax tree

**IR generator**

intermediate representation

```
tmp := ADD a a
    := SUB tmp zo
```

**optimiser**

intermediate representation

```
tmp := SHL a 1
fu := SUB tmp zo
```

**code generator**

assembly code

```
sll $8, $4, 1
sub $2, $8, $5
```

# Live-variable analysis

```
int main (int a, int b, int c) {
  int d, e, f, g, h;

  d = a + b;

  e = d + b;

  f = d + a;

  g = a + c;

  h = e + g;

  return f + h;

}
```

# Live-variable analysis

```
int main (int a, int b, int c) {
  int d, e, f, g, h;
  //{ a, b, c }
  d = a + b;

  e = d + b;

  f = d + a;

  g = a + c;

  h = e + g;

  return f + h;

}
```

# Live-variable analysis

```
int main (int a, int b, int c) {
  int d, e, f, g, h;
  //{ a, b, c }
  d = a + b;
  //{ a, b, c, d }
  e = d + b;


  f = d + a;


  g = a + c;


  h = e + g;


  return f + h;


}
```

# Live-variable analysis

```
int main (int a, int b, int c) {
  int d, e, f, g, h;
  //{ a, b, c }
  d = a + b;
  //{ a, b, c, d }
  e = d + b;
  //{ a, c, d, e }
  f = d + a;


  g = a + c;


  h = e + g;


  return f + h;


}
```

# Live-variable analysis

```
int main (int a, int b, int c) {
  int d, e, f, g, h;
  //{ a, b, c }
  d = a + b;
  //{ a, b, c, d }
  e = d + b;
  //{ a, c, d, e }
  f = d + a;
  //{ a, c, e, f }
  g = a + c;

  h = e + g;

  return f + h;

}
```

# Live-variable analysis

```
int main (int a, int b, int c) {
  int d, e, f, g, h;
  //{ a, b, c }
  d = a + b;
  //{ a, b, c, d }
  e = d + b;
  //{ a, c, d, e }
  f = d + a;
  //{ a, c, e, f }
  g = a + c;
  //{ e, f, g }
  h = e + g;

  return f + h;

}
```

# Live-variable analysis

```
int main (int a, int b, int c) {
  int d, e, f, g, h;
  //{ a, b, c }
  d = a + b;
  //{ a, b, c, d }
  e = d + b;
  //{ a, c, d, e }
  f = d + a;
  //{ a, c, e, f }
  g = a + c;
  //{ e, f, g }
  h = e + g;
  //{ f, h }
  return f + h;

}
```

# Live-variable analysis

```
int main (int a, int b, int c) {
  int d, e, f, g, h;
  //{ a, b, c }
  d = a + b;
  //{ a, b, c, d }
  e = d + b;
  //{ a, c, d, e }
  f = d + a;
  //{ a, c, e, f }
  g = a + c;
  //{ e, f, g }
  h = e + g;
  //{ f, h }
  return f + h;
  //∅
}
```

# Live-variable analysis

```
int main (int a, int b, int c) {
  int d, e, f, g, h;

  d = a + b;

  e = d + b;

  f = d + a;

  g = a + c;

  h = e + g;

  return f + h;

}
```

# Live-variable analysis

```
int main (int a, int b, int c) {
  int d, e, f, g, h;

  d = a + b;

  e = d + b;

  f = d + a;

  g = a + c;

  h = e + g;

  return f + h;
  // ∅
}
```

# Live-variable analysis

```
int main (int a, int b, int c) {
  int d, e, f, g, h;

  d = a + b;

  e = d + b;

  f = d + a;

  g = a + c;

  h = e + g;
  //{ f, h }
  return f + h;
  // ∅
}
```

# Live-variable analysis

```
int main (int a, int b, int c) {
  int d, e, f, g, h;

  d = a + b;

  e = d + b;

  f = d + a;

  g = a + c;
  //{ e, f, g }
  h = e + g;
  //{ f, h }
  return f + h;
  // ∅
}
```

# Live-variable analysis

```
int main (int a, int b, int c) {
  int d, e, f, g, h;

  d = a + b;

  e = d + b;

  f = d + a;
  //{ a, c, e, f }
  g = a + c;
  //{ e, f, g }
  h = e + g;
  //{ f, h }
  return f + h;
  //∅
}
```

# Live-variable analysis

```
int main (int a, int b, int c) {
  int d, e, f, g, h;

  d = a + b;

  e = d + b;
  //{ a, c, d, e }
  f = d + a;
  //{ a, c, e, f }
  g = a + c;
  //{ e, f, g }
  h = e + g;
  //{ f, h }
  return f + h;
  //∅
}
```

# Live-variable analysis

```
int main (int a, int b, int c) {
  int d, e, f, g, h;

  d = a + b;
  //{ a, b, c, d }
  e = d + b;
  //{ a, c, d, e }
  f = d + a;
  //{ a, c, e, f }
  g = a + c;
  //{ e, f, g }
  h = e + g;
  //{ f, h }
  return f + h;
  //∅
}
```

# Live-variable analysis

```
int main (int a, int b, int c) {
  int d, e, f, g, h;
  //{ a, b, c }
  d = a + b;
  //{ a, b, c, d }
  e = d + b;
  //{ a, c, d, e }
  f = d + a;
  //{ a, c, e, f }
  g = a + c;
  //{ e, f, g }
  h = e + g;
  //{ f, h }
  return f + h;
  //∅
}
```

$$\text{live}_{before}(s) = \text{live}_{after}(s) - \text{def}_V(s) \cup \text{use}_V(s)$$

$$\text{def}_V(\textbf{return } \texttt{f + h}) = \emptyset$$
$$\text{def}_V(\texttt{h = e + g}) \quad\ = \{\, h \,\}$$
$$\text{use}_V(\textbf{return } \texttt{f + h}) = \{\, f, h \,\}$$
$$\text{use}_V(\texttt{h = e + g}) \quad\ = \{\, e, g \,\}$$

# Live-variable analysis

- Useful for **register allocation**.

- Also useful for **dead code elimination**.

```
int main (int a, int b, int c) {
  int d, e;

  d = a + b;

  e = b * c;

  return a + e;

}
```

# Live-variable analysis

- Useful for **register allocation**.

- Also useful for **dead code elimination**.

```
int main (int a, int b, int c) {
  int d, e;

  d = a + b;

  e = b * c;

  return a + e;
  // ∅
}
```

# Live-variable analysis

- Useful for **register allocation**.

- Also useful for **dead code elimination**.

```
int main (int a, int b, int c) {
  int d, e;

  d = a + b;


  e = b * c;
  //{a,e}
  return a + e;
  //∅
}
```

# Live-variable analysis

- Useful for **register allocation**.

- Also useful for **dead code elimination**.

```
int main (int a, int b, int c) {
  int d, e;

  d = a + b;
  //{a, b, c}
  e = b * c;
  //{a, e}
  return a + e;
  //∅
}
```

# Live-variable analysis

- Useful for **register allocation**.

- Also useful for **dead code elimination**.

```
int main (int a, int b, int c) {
  int d, e;
  //{ a, b, c }
  d = a + b;
  //{ a, b, c }
  e = b * c;
  //{ a, e }
  return a + e;
  //∅
}
```

| Analysis | Operates on | Transformation | | |
|---|---|---|---|---|
| live variables | sets of variables | register allocation, dead code elimination | | |
| | | | | |
| | | | | |
| | | | | |

# Available expressions

```
int main (int a, int b, int c) {
  int d, e, f;

  d = a + b;

  e = c * b;

  f = a + b;

  a = f / 10;

  return a + b;
}
```

# Available expressions

```
int main (int a, int b, int c) {
  int d, e, f;
  // ∅
  d = a + b;

  e = c * b;

  f = a + b;

  a = f / 10;

  return a + b;
}
```

# Available expressions

```
int main (int a, int b, int c) {
  int d, e, f;
  // ∅
  d = a + b;
  // { a+b }
  e = c * b;

  f = a + b;


  a = f / 10;


  return a + b;
}
```

# Available expressions

```
int main (int a, int b, int c) {
  int d, e, f;
  // ∅
  d = a + b;
  // { a+b }
  e = c * b;
  // { a+b, c*b }
  f = a + b;

  a = f / 10;

  return a + b;
}
```
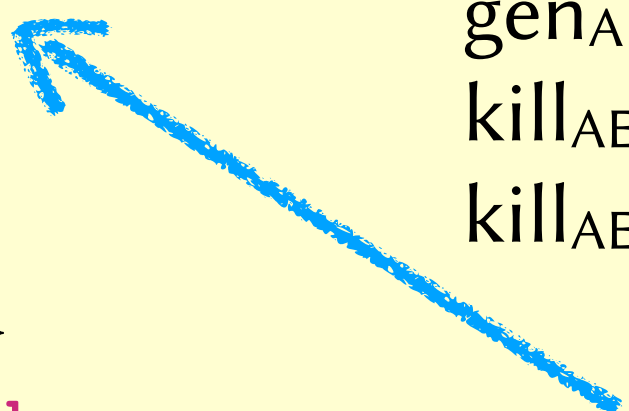
# Available expressions

```
int main (int a, int b, int c) {
  int d, e, f;
  // ∅
  d = a + b;
  // { a+b }
  e = c * b;
  // { a+b, c*b }
  f = a + b;
  // { a+b, c*b }
  a = f / 10;

  return a + b;
}
```

# Available expressions

```
int main (int a, int b, int c) {
  int d, e, f;
  // ∅
  d = a + b;
  // { a+b }
  e = c * b;
  // { a+b, c*b }
  f = a + b;
  // { a+b, c*b }
  a = f / 10;
  // { c*b, f/10 }
  return a + b;
}
```

$$\text{avail}_{\text{after}}(s) = \text{avail}_{\text{before}}(s) \cup \text{gen}_{\text{AE}}(s) - \text{kill}_{\text{AE}}(s)$$

$\text{gen}_{\text{AE}}(\textbf{return}\ a\ +\ b) = \{\ a+b\ \}$
$\text{gen}_{\text{AE}}(a\ =\ f\ /\ 10)\quad = \{\ f/10\ \}$
$\text{kill}_{\text{AE}}(\textbf{return}\ a\ +\ b)\ = \emptyset$
$\text{kill}_{\text{AE}}(a\ =\ f\ /\ 10)\quad = E_a$

# Available expressions

```
int main (int a, int b, int c) {
  int d, e, f;
  // ∅
  d = a + b;
  // { a+b }
  e = c * b;
  // { a+b, c*b }
  f = d;
  // { a+b, c*b }
  a = f / 10;
  // { c*b, f/10 }
  return a + b;
}
```

$$\text{avail}_{after}(s) = \text{avail}_{before}(s) \cup \text{gen}_{AE}(s) - \text{kill}_{AE}(s)$$

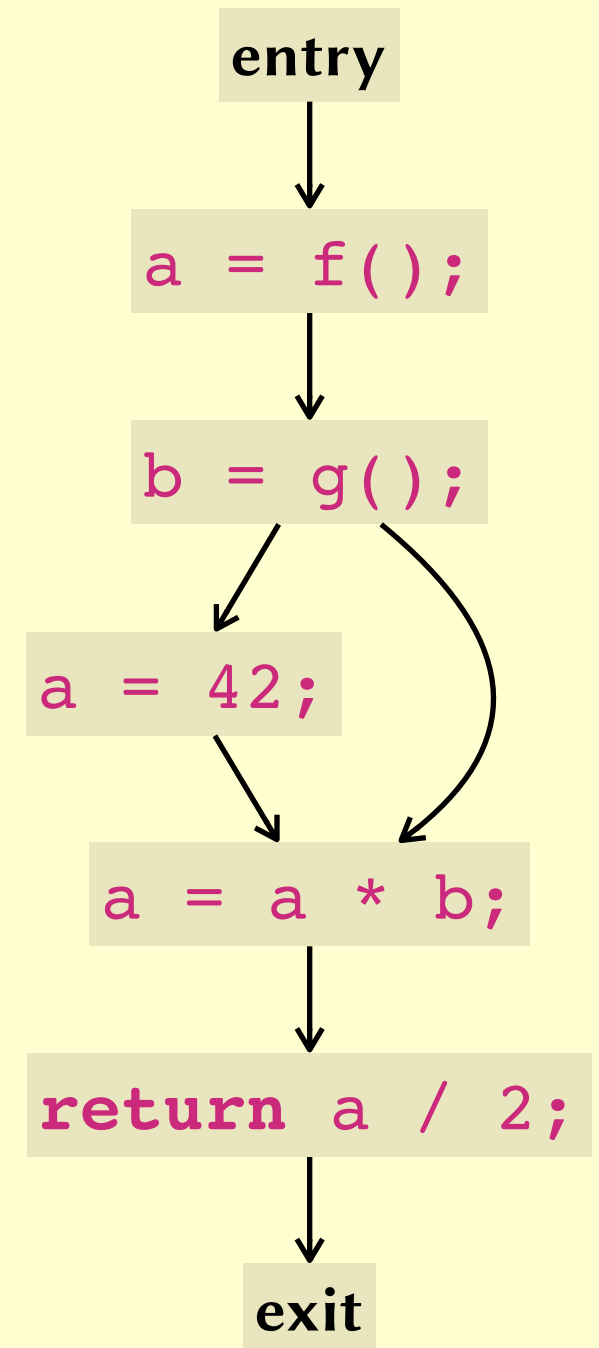$\text{gen}_{AE}(\textbf{return } a + b) = \{ a+b \}$
$\text{gen}_{AE}(a = f / 10) = \{ f/10 \}$
$\text{kill}_{AE}(\textbf{return } a + b) = \varnothing$
$\text{kill}_{AE}(a = f / 10) = E_a$

"common subexpression elimination"

| Analysis | Operates on | Transformation | Direction | |
|---|---|---|---|---|
| live variables | sets of variables | register allocation, dead code elimination | Backward | |
| available expressions | sets of expressions | common subexpression elimination | Forward | |
| | | | | |
| | | | | |

# Control-flow graphs

```
a = f();
b = g();
if (b < 0) {
  a = 42;
}
a = a * b;
return a / 2;
```

➡️

**entry**

a = f();

b = g();

b<0     b≥0

a = 42;

a = a * b;

**return** a / 2;

**exit**

➡️

**entry**

a = f();

b = g();

a = 42;

a = a * b;

**return** a / 2;

**exit**

# Control-flow graphs



```
entry

a = f();

b = g();
  b<0        b≥0

a = 42;

a = a * b;

return a / 2;

exit
```

```
entry

a = f();

b = g();

a = 42;

a = a * b;

return a / 2;

exit
```

"basic block"

```
entry

a = f();
b = g();

a = 42;

a = a * b;
return a / 2;

exit
```

# Terminology

entry

B0:
```
a = f();
b = g();
```

B1:
```
a = 42;
```

B2:
```
a = a * b;
return a / 2;
```

exit

succ(B) = all blocks that B can jump to

pred(B) = all blocks that can jump to B

analyseBlock(B, X) = analyse block B, starting from X

# LV on a CFG

*in:*
**entry**
*out:*

*in:*
B0: `a = f();`
`b = g();`
*out:* a  b

*in:* b
B1: `a = 42;`
*out:* a  b

*in:* a  b
B2: `a = a * b;`
`return a / 2;`
*out:*

*in:*
**exit**
*out:*

$$out[B] = in[B_1] \cup \dots \cup in[B_n]$$
$$\text{where } succ(B) = \{B_1, \dots, B_n\}$$

$$in[B] = \text{analyseBlock}_{LV}(B, out[B])$$

Start with $in[B] = out[B] = \varnothing$, for all B. Then keep applying these definitions until nothing changes.

*Gary Kildall*
*1942–1994*

# Analysis precision

```
int main(int a,
         int b) {
  int foo[2];
  foo[a] = 42;
  return foo[b];
}
```

```
int main() {
  int a = 42;
  if (1) return 0;
  else return a;
}
```

```
int main() {
  int a = 42;
  while(1);
  return a;
}
```

- What happens if the analysis wrongly says a variable **is not live**?

- What happens if the analysis wrongly says that a variable **is live**?

- Live-variable analysis **overapproximates** the set of live variables.

# AE on a CFG

*in:*
**entry**
*out:*

BB1:
*in:*
```
c = a + b;
f = b * c;
```
*out:* a+b  b*c  a*b

BB2:
*in:*
```
d = a * b;
print(a + b);
```
*out:* a+b  a*b  b*c

BB3:
*in:* a+b
```
e = a * b;
```
*out:* a*b  a+b  b*c

BB4:
*in:* a*b  a+b
```
f = a + b;
```
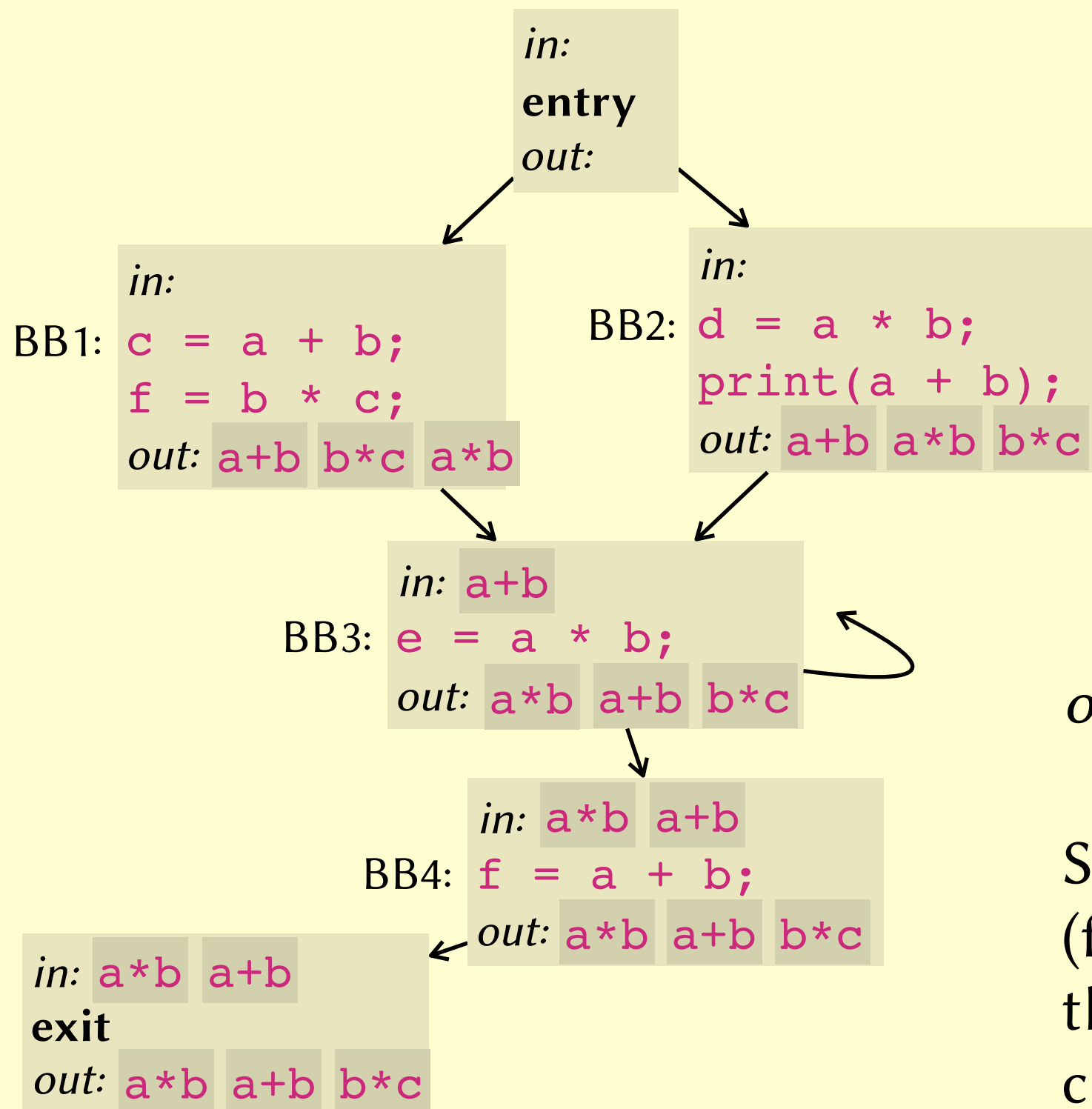*out:* a*b  a+b  b*c

*in:* a*b  a+b
**exit**
*out:* a*b  a+b  b*c

$in[B] = out[B_1] \cap \ldots \cap out[B_n]$
where pred(B) = {$B_1,\ldots,B_n$}

$out[B] = \text{analyseBlock}_{AE}(in[B])$

Start with *out*[B] = *all expressions* (for all B). Then keep applying these definitions until nothing changes.

# AE on a CFG

*in:*
**entry**
*out:*

BB1:
*in:*
```
c = a + b;
f = b * c;
```
*out:* a+b b*c

BB2:
*in:*
```
d = a * b;
print(a + b);
```
*out:* a+b a*b

BB3:
*in:* a+b
```
e = a * b;
```
*out:* a*b a+b

BB4:
*in:* a*b a+b
```
f = a + b;
```
*out:* a*b a+b

*in:* a*b a+b
**exit**
*out:* a*b a+b

$$in[B] = out[B_1] \cap \ldots \cap out[B_n]$$
$$\text{where pred}(B) = \{B_1, \ldots, B_n\}$$

$$out[B] = \text{analyseBlock}_{AE}(in[B])$$

Start with *out*[B] = *all expressions* (for all B). Then keep applying these definitions until nothing changes.

# AE on a CFG

*in:*
**entry**
*out:*

*in:*
BB0: `t = a + b;`
*out:* `a+b`

*in:* `a+b`
BB1: `c = a + b;`
`f = b * c;`
*out:* `a+b` `b*c`

*in:* `a+b`
BB2: `d = a * b;`
`print(a + b);`
*out:* `a+b` `a*b`

*in:* `a+b`
BB3: `e = a * b;`
*out:* `a*b` `a+b`

*in:* `a*b` `a+b`
BB4: `f = a + b;`
*out:* `a*b` `a+b`

*in:* `a*b` `a+b`
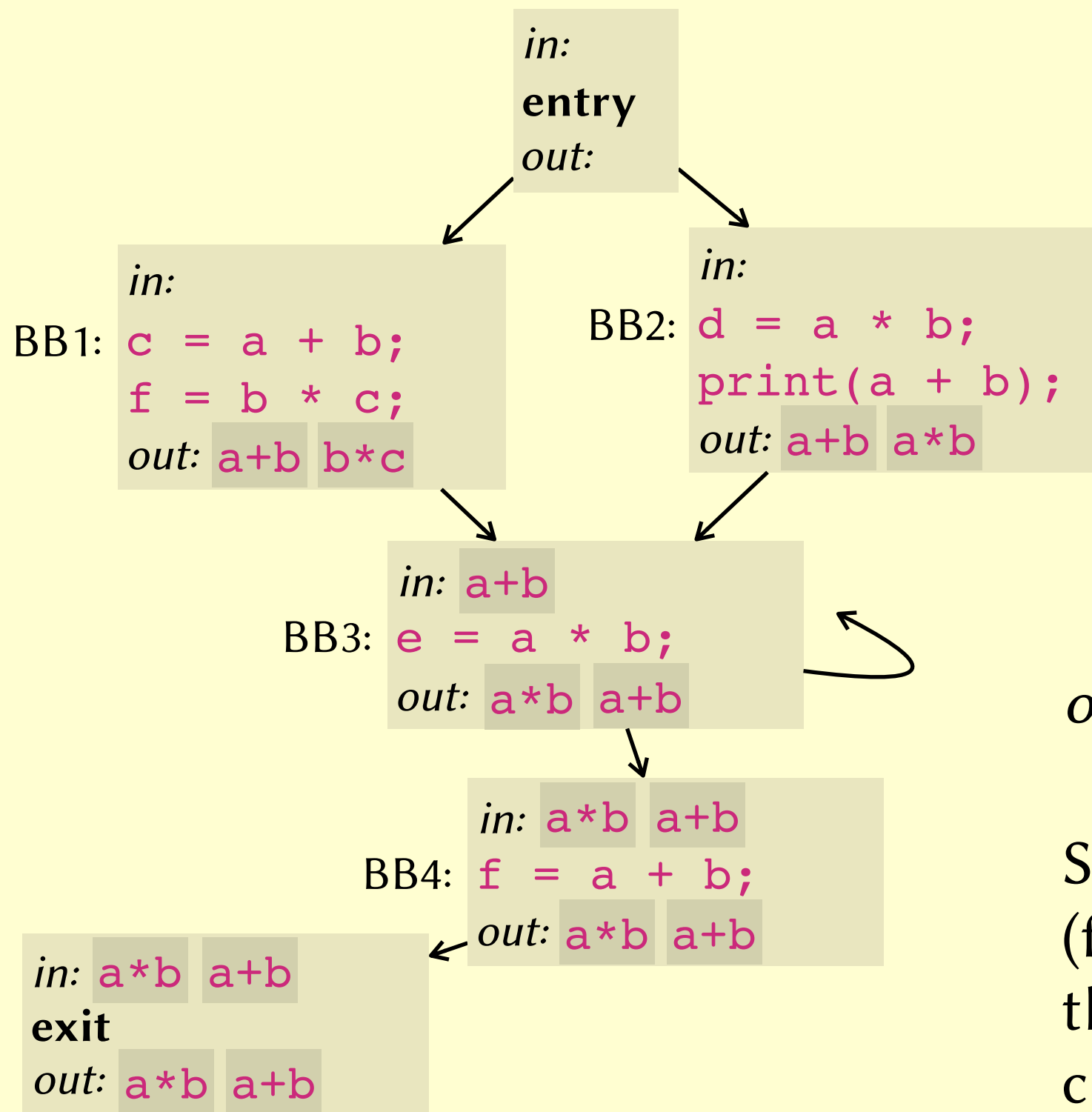**exit**
*out:* `a*b` `a+b`

$$in[B] = out[B_1] \cap \ldots \cap out[B_n]$$
$$\text{where pred}(B) = \{B_1, \ldots, B_n\}$$

$$out[B] = \text{analyseBlock}_{AE}(in[B])$$

Start with *out*[B] = *all expressions* (for all B). Then keep applying these definitions until nothing changes.
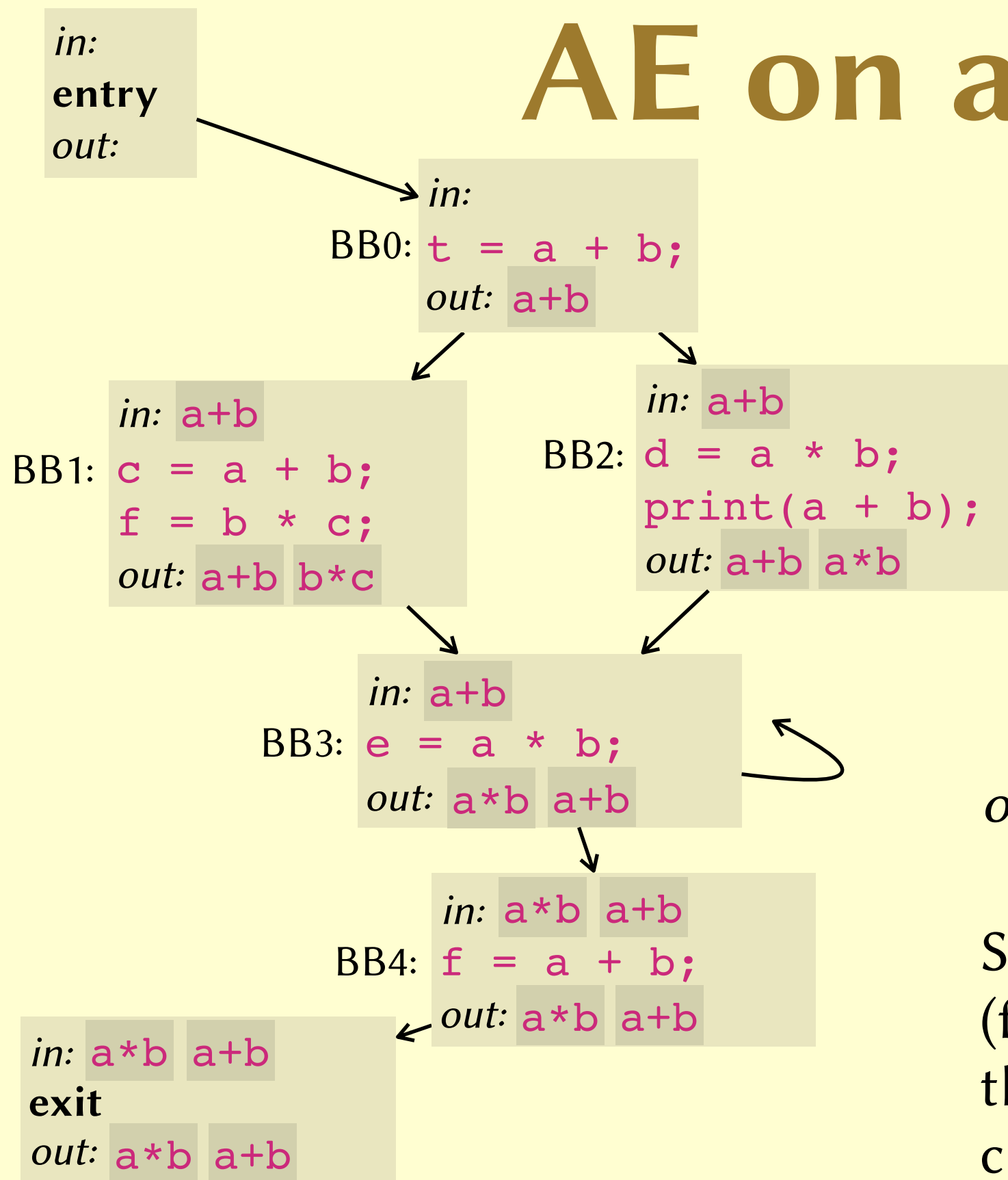
# AE on a CFG

*in:*
**entry**
*out:*

*in:*
BB0: `t = a + b;`
*out:* a+b

*but was it worth it?*

*in:* a+b
BB1: `c = t;`
`f = b * c;`
*out:* a+b  b*c

*in:* a+b
BB2: `d = a * b;`
`print(t);`
*out:* a+b  a*b

*in:* a+b
BB3: `e = a * b;`
*out:* a*b  a+b

*in:* a*b  a+b
BB4: `f = t;`
*out:* a*b  a+b

*in:* a*b  a+b
**exit**
*out:* a*b  a+b

$in[B] = out[B_1] \cap \dots \cap out[B_n]$
$\quad\quad\quad\quad$ where $pred(B) = \{B_1, \dots, B_n\}$

$out[B] = \text{analyseBlock}_{AE}(in[B])$

Start with *out*[B] = *all expressions* (for all B). Then keep applying these definitions until nothing changes.

# Analysis precision
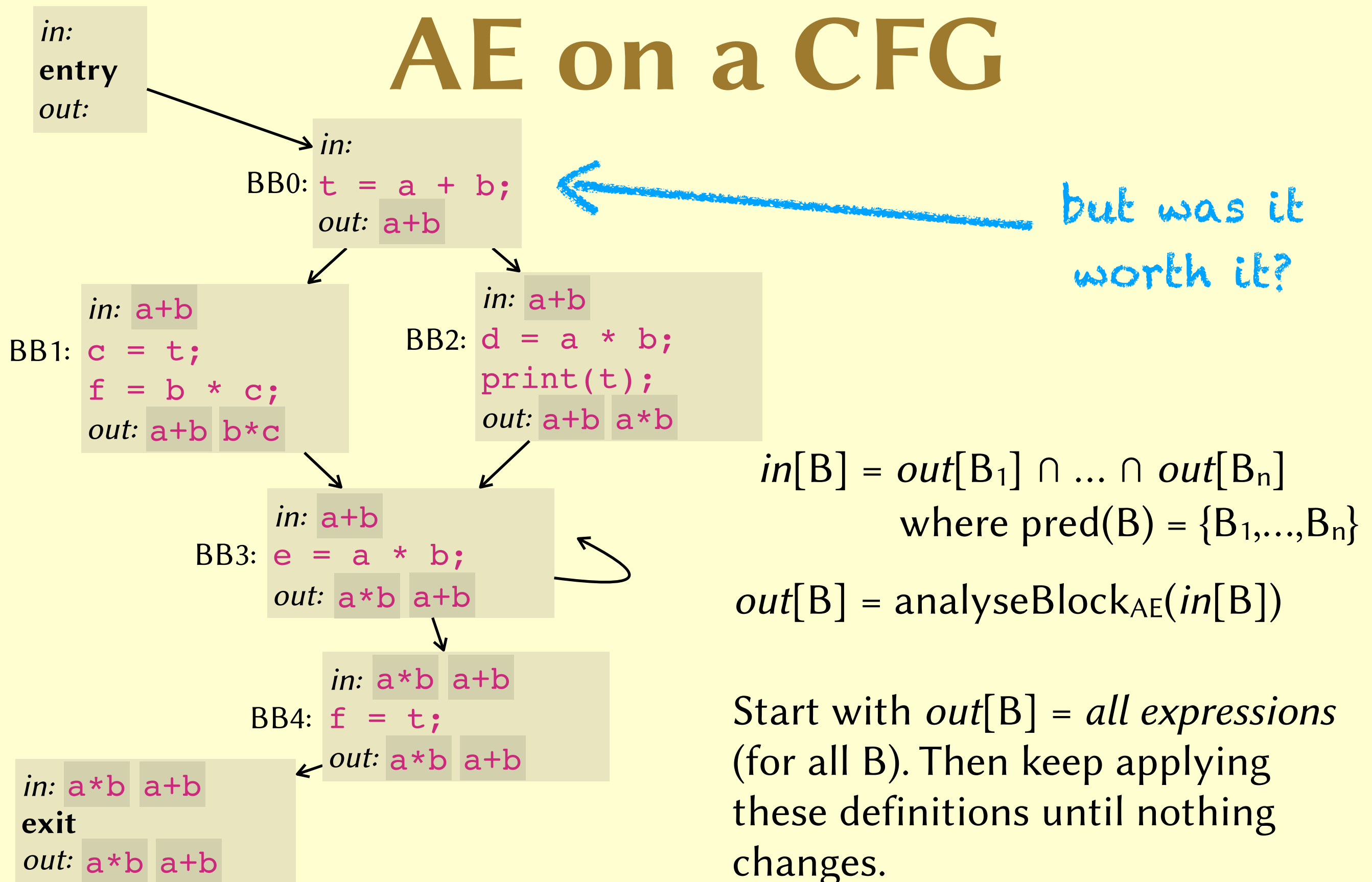
- What happens if the analysis wrongly says an expression **is not available**?

- What happens if the analysis wrongly says that an expression **is available**?

- Available-expression analysis **underapproximates** the set of available expressions.

| Analysis | Operates on | Transformation | Direction | Meet |
|---|---|---|---|---|
| live variables | sets of variables | register allocation, dead code elimination | Backward | $\cup$ |
| available expressions | sets of expressions | common subexpression elimination | Forward | $\cap$ |
|  |  |  |  |  |
|  |  |  |  |  |

# Reaching definitions

```
23  int main () {

24    int a = f();

25    int b = g();

26    if (b > 0) a = 42;

27    a = a * b;

28    return a / 2;
    }
```

# Reaching definitions

```
23 int main () {
   // ∅
24    int a = f();

25    int b = g();

26    if (b > 0) a = 42;

27    a = a * b;

28    return a / 2;
   }
```

# Reaching definitions

```
23  int main () {
      // ∅
24    int a = f();
      // { 24: a }
25    int b = g();

26    if (b > 0) a = 42;

27    a = a * b;

28    return a / 2;
    }
```

# Reaching definitions

```
23  int main () {
      // ∅
24    int a = f();
      // { 24: a }
25    int b = g();
      // { 24: a, 25: b }
26    if (b > 0) a = 42;


27    a = a * b;


28    return a / 2;
    }
```
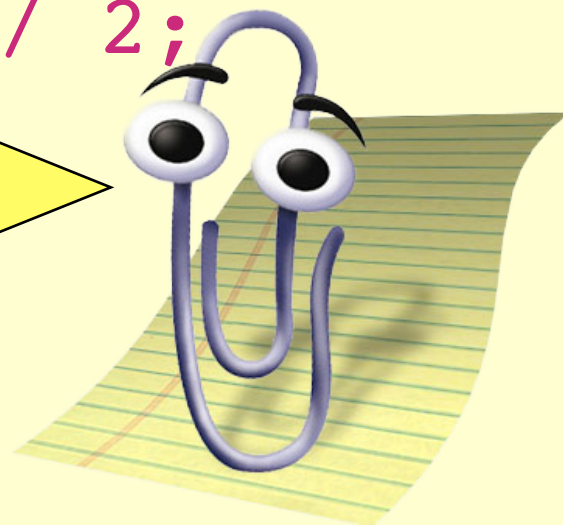
# Reaching definitions

```
23  int main () {
      // ∅
24    int a = f();
      // { 24: a }
25    int b = g();
      // { 24: a, 25: b }
26    if (b > 0) a = 42;
      // { 24: a, 25: b, 26: a }
27    a = a * b;

28    return a / 2;
    }
```

# Reaching definitions

```
23  int main () {
        // ∅
24      int a = f();
        // { 24: a }
25      int b = g();
        // { 24: a, 25: b }
26      if (b > 0) a = 42;
        // { 24: a, 25: b, 26: a }
27      a = a * b;
        // { 25: b, 27: a }
28      return a / 2;
    }
```

```
int main () {
    int a = 7;
    int b = f();
    return a - b;
}
```

**The value of a on line 28 is provided on line 27**

# Reaching definitions

```
23  int main () {
      // ∅
24    int a = f();
      // { 24: a }
25    int b = g();
      // { 24: a, 25: b }
26    if (b > 0) a = 42;
      // { 24: a, 25: b, 26: a }
27    a = a * b;
      // { 25: b, 27: a }
28    return a / 2;
    }
```

```
int main () {
    int a = 7;
    int b = f();
    return 7 - b;
}
```

"constant propagation"

The value of **a** on line 28 is provided on line 27

Oops, you forgot to initialise **d**!

```
int main () {
    int a = f();
    return a - d;
}
```

# Reaching definitions

```
23  int main () {
     //∅
24    int a = f();
     //{ 24: a }
25    int b = g();
     //{ 24: a, 25: b }
26    if (b > 0) a = 42;
     //{ 24: a, 25: b, 26: a }
27    a = a * b;
     //{ 25: b, 27: a }
28    return a / 2;
    }
```

reach$_{after}$(s) = reach$_{before}$(s) - kill$_D$(s) ∪ gen$_D$(s)

kill$_D$(*L:* **return** a / 2) = ∅
kill$_D$(*L:* a = 42)           = D$_a$
gen$_D$(*L:* **return** a / 2) = ∅
gen$_D$(*L:* a = 42)           = { *L:* a }

# Reaching definitions

*in:*
**entry**
*out:*

*in:*
BB0:  *24* `a = f();`
      *25* `b = g();`
      *out:* a:*24*  b:*25*

*in:* a:*24*  b:*25*
BB1:  *26* `a = 42;`
      *out:* a:*26*  b:*25*

*in:* a:*24*  a:*26*  b:*25*
BB2:  *27* `a = a * b;`
      *28* **return** `a / 2;`
      *out:* a:*27*  b:*25*

*in:* a:*27*  b:*25*
**exit**
*out:* a:*27*  b:*25*

$in[B] = out[B_1] \cup ... \cup out[B_n]$
   where $pred(B) = \{B_1,...,B_n\}$

$out[B] = analyseBlock_{RD}(in[B])$

Start with $out[B] = in[B] = \emptyset$, for all B. Then keep applying these definitions until nothing changes.

# Dataflow analyses

Frances Allen
1932-2020

| Analysis | Operates on | Transformation | Direction | Meet |
|---|---|---|---|---|
| live variables | sets of variables | register allocation, dead code elimination | Backward | ∪ |
| available expressions | sets of expressions | common subexpression elimination | Forward | ∩ |
| reaching definitions | sets of definitions | constant propagation, spotting undefined variables | Forward | ∪ |
| | | | | |

# Dataflow analyses

Frances Allen
1932-2020

| Analysis | Operates on | Transformation | Direction | Meet |
|---|---|---|---|---|
| live variables | sets of variables | register allocation, dead code elimination | Backward | ∪ |
| available expressions | sets of expressions | common subexpression elimination | Forward | ∩ |
| reaching definitions | sets of definitions | constant propagation, spotting undefined variables | Forward | ∪ |
| ? | ? | ? | Backward | ∩ |

# Very busy expressions

```
int main (int a, int b, int c) {

  if (b > 0) {
    print(b + c);
  } else if (b < 0) {
    f(b + c);
  } else {
    a = b + c;
  }
  a = a * a;
  return a;
}
```

# Very busy expressions

```
int main (int a, int b, int c) {
  int t = b + c;
  if (b > 0) {
    print(t);
  } else if (b < 0) {
    f(t);
  } else {
    a = t;
  }
  a = a * a;
  return a;
}
```

"hoisting"

# Summary

- Data flow analysis approximates a program by a **flow graph** made up of **basic blocks**.

- Analyses can go **forwards** or **backwards**.

- Analyses can consider **all** ($\cap$) paths or **any** ($\cup$) paths.

- Analyses may not be completely **precise**, so need to **overapproximate** or **underapproximate** as appropriate.

- Analyses enable **transformations** (but need to consider whether each transformation is actually worth doing).