

Compiler Design

Lecture 19: Instruction Selection via Tree-pattern matching

Christophe Dubach
Winter 2025

(EaC-11.3)

Copyright 2003, Keith D. Cooper, Ken Kennedy & Linda Torczon, all rights reserved.

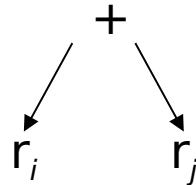
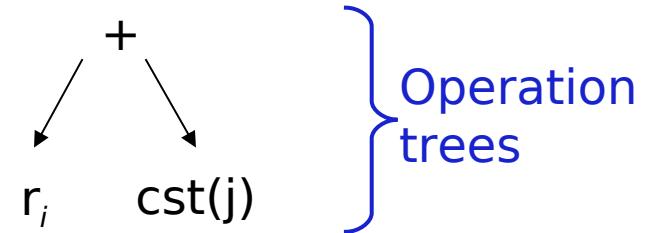
The Concept

Many compilers use tree-structured IRs

- Abstract syntax trees generated in the parser
- Trees or DAGs for expressions

These systems might well use trees to represent target ISA

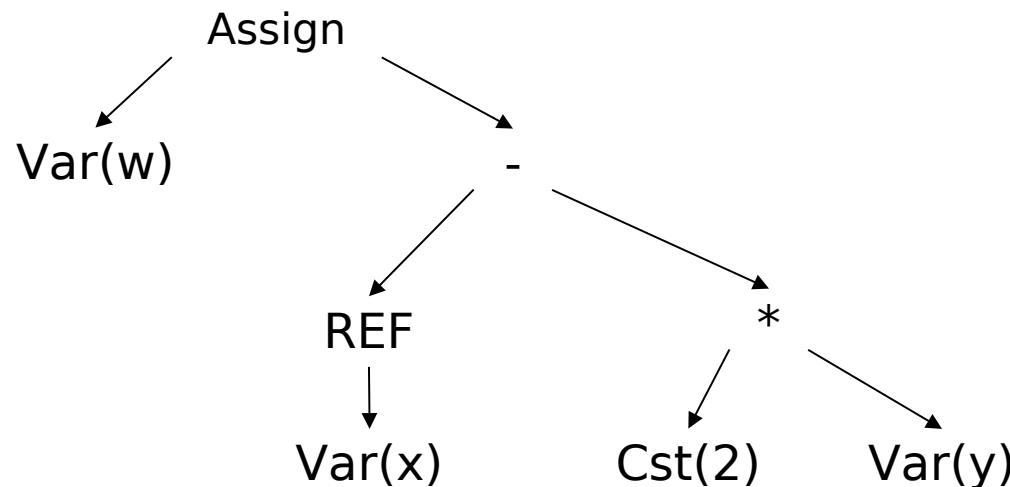
Consider the add operators

$$\text{add } r_i, r_j \Rightarrow r_k$$

$$\text{addI } r_i, j \Rightarrow r_k$$


What if we could match these “operation trees” against IR tree?

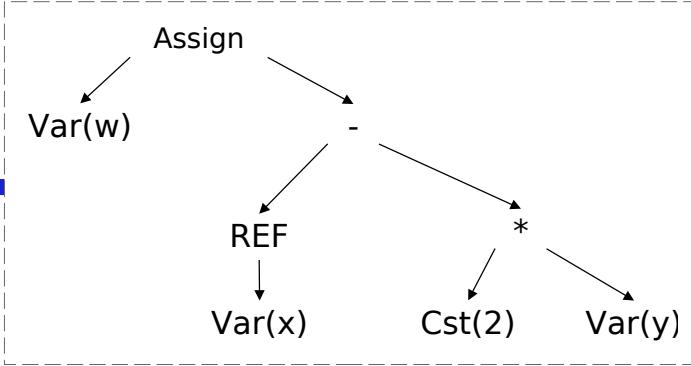
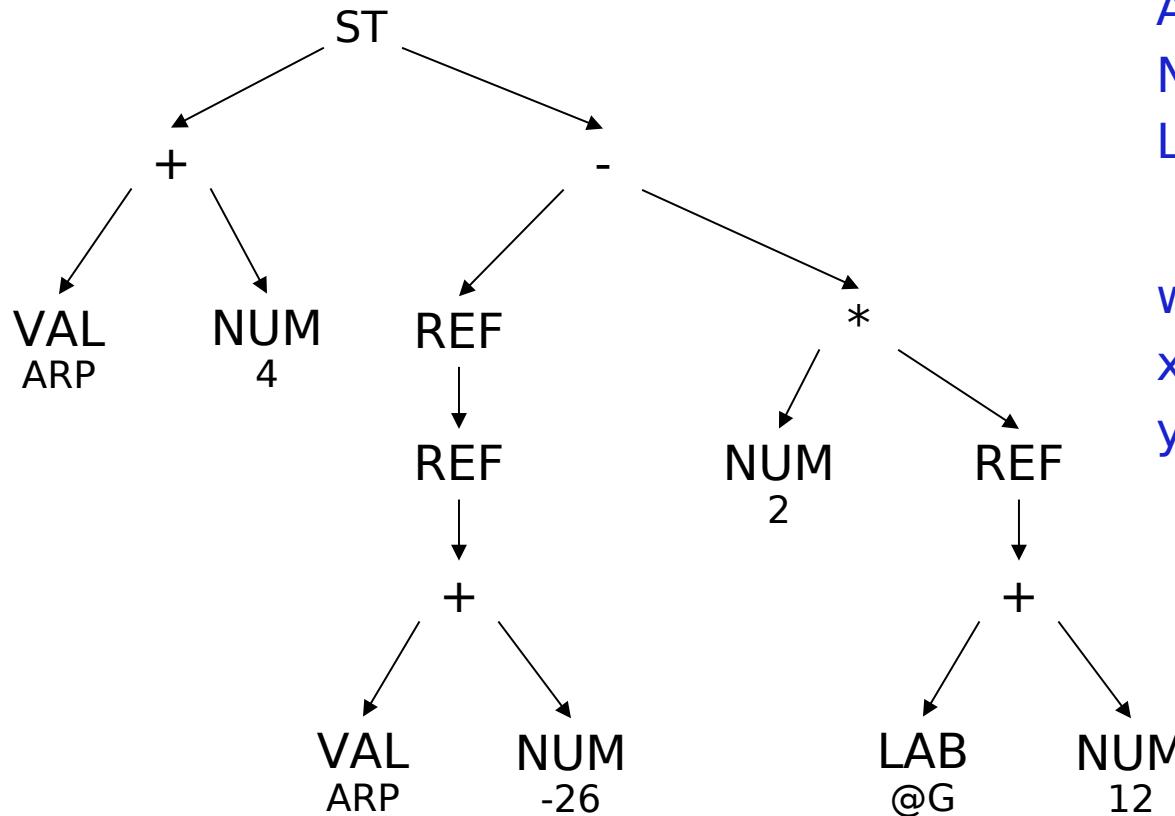
The Concept

AST for $w \leftarrow (*x) - 2 * y$



The Concept

Low-level AST for $w \leftarrow (*x) - 2 * y$



ARP: \$fp

NUM: constant

LAB: ASM label

w: at ARP+4

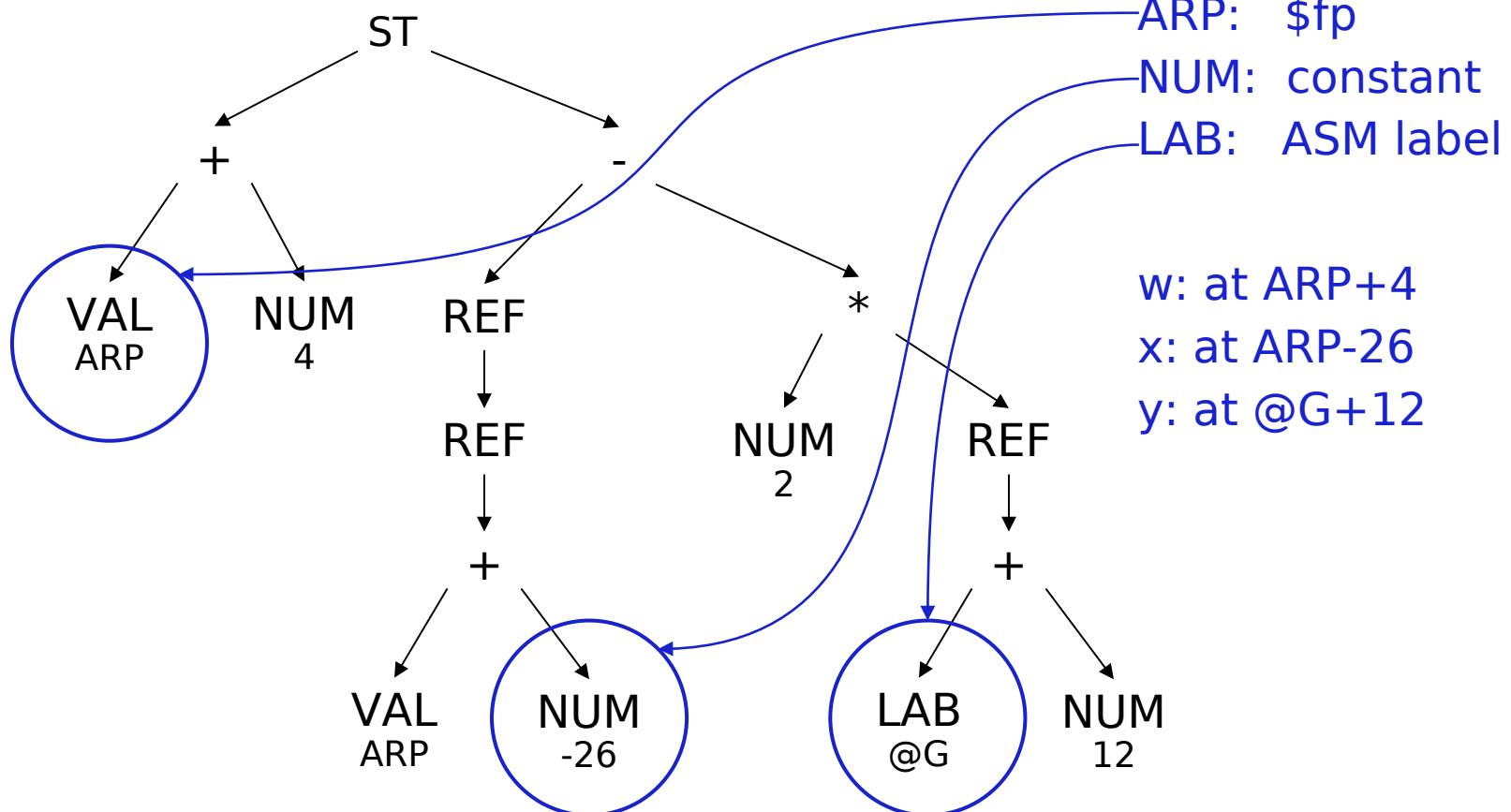
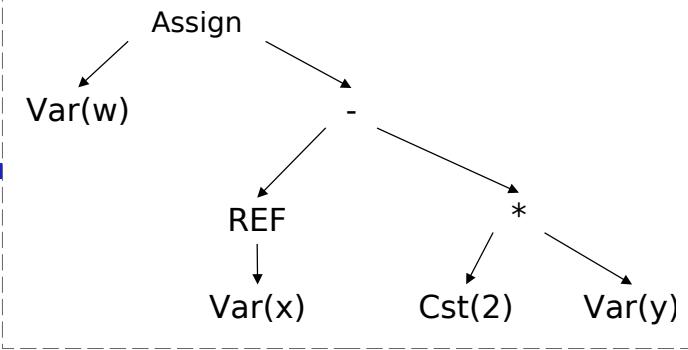
x: at ARP-26

y: at @G+12

ARP = Activation Record Pointer = frame pointer

The Concept

Low-level AST for $w \leftarrow (*x) - 2 * y$



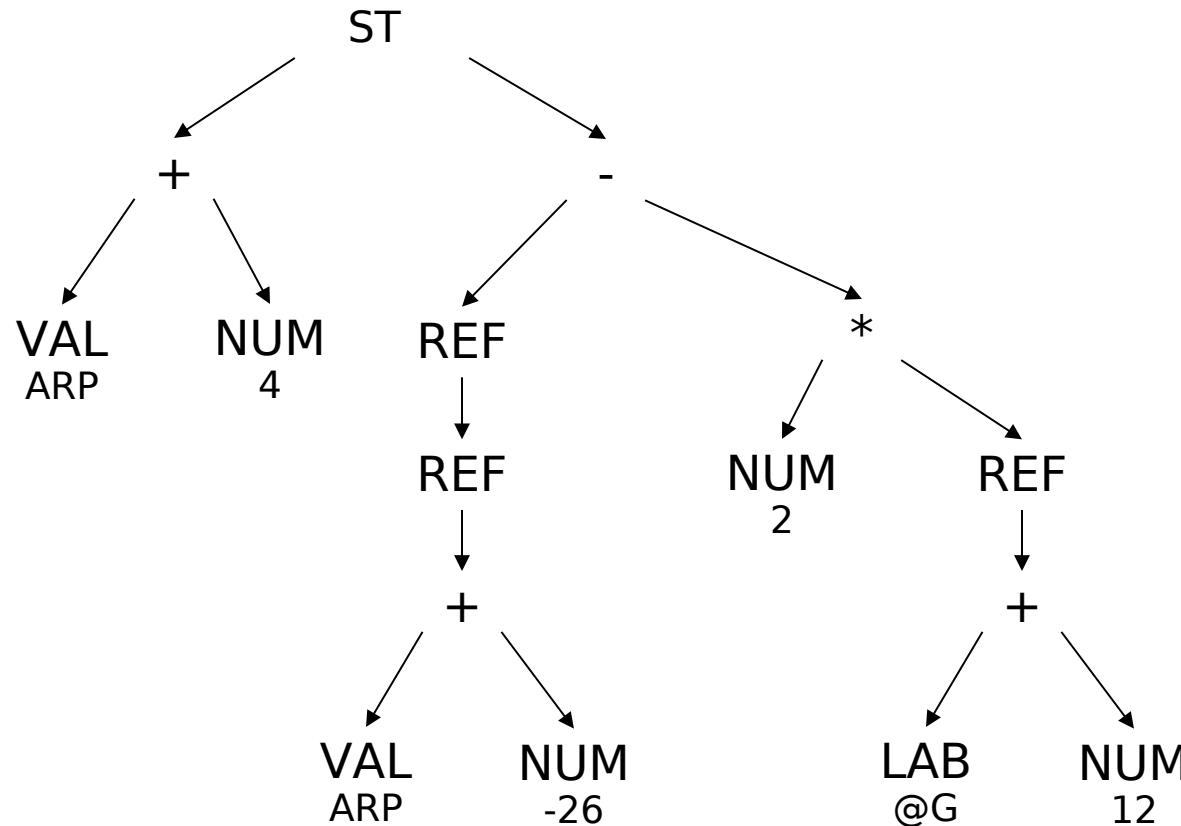
ARP = Activation Record Pointer = frame pointer

Tree-pattern matching

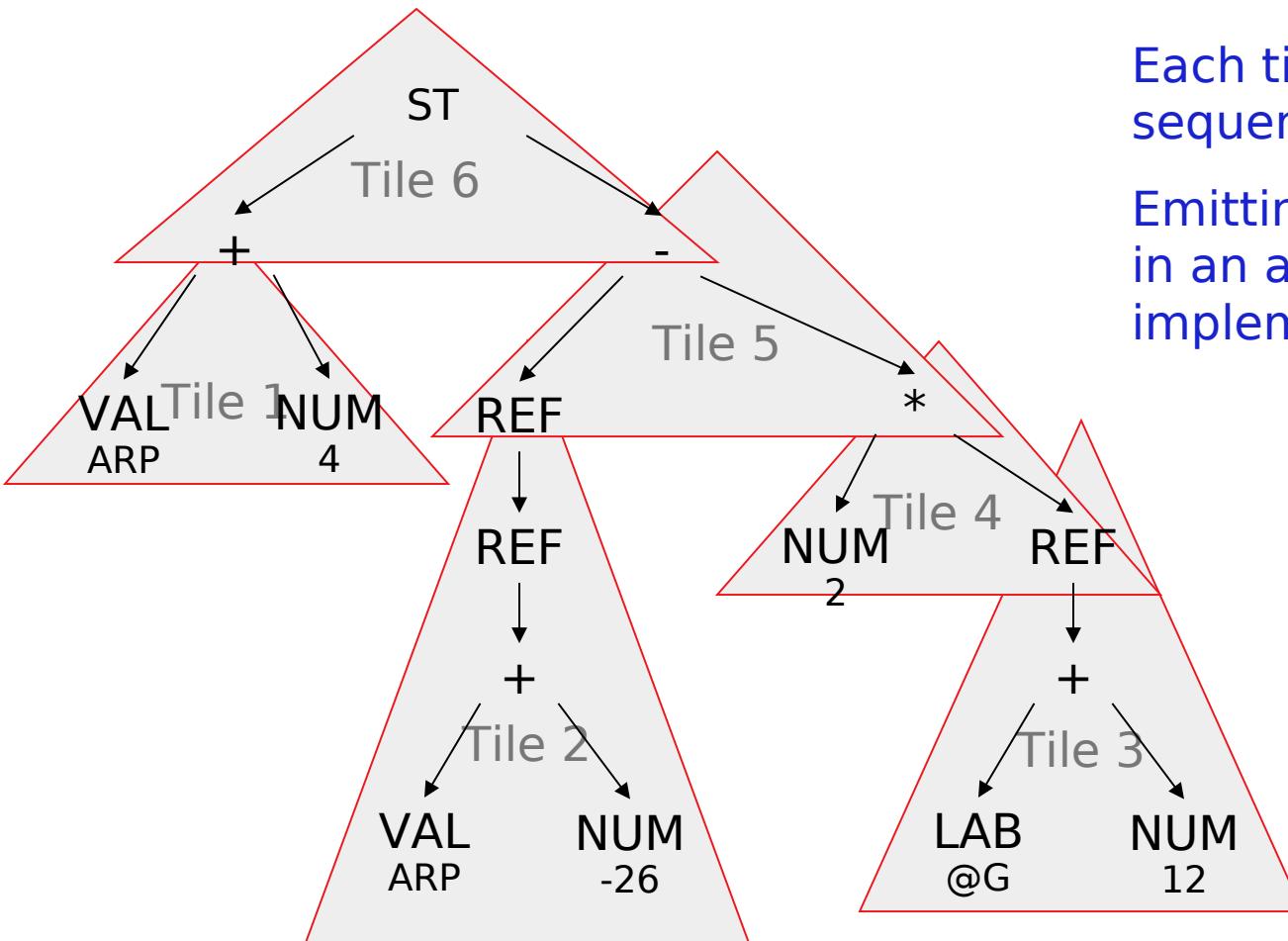
Goal is to “tile” AST with operation trees

- A tiling is collection of $\langle \text{ast}, \text{op} \rangle$ pairs
 - ast is a node in the low-level AST
 - op is an operation tree
 - $\langle \text{ast}, \text{op} \rangle$ means that op could implement the subtree at ast
- A tiling ‘implements’ an AST if it covers every node in the AST and the overlap between any two trees is limited to a single node
 - $\langle \text{ast}, \text{op} \rangle \in \text{tiling}$ means ast is also covered by a leaf in another operation tree in the tiling, unless it is the root
 - Where two operation trees meet, they must be compatible (expect the value in the same location)

Tiling the Tree



Tiling the Tree

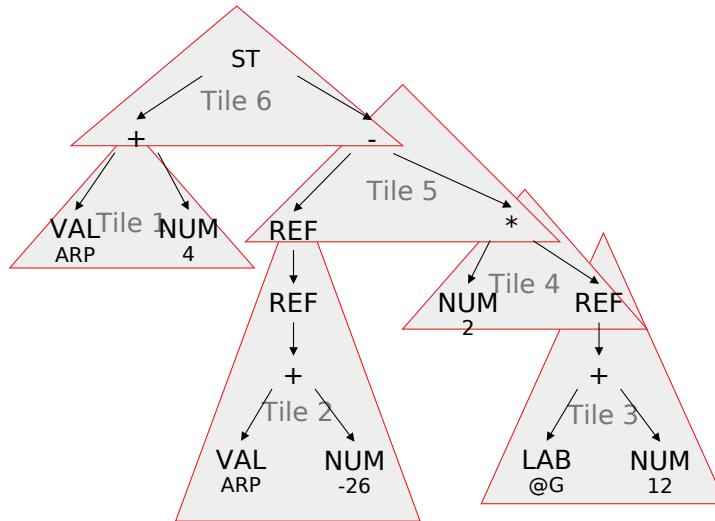


Each tile corresponds to a sequence of operations
Emitting those operations in an appropriate order
implements the tree.

Generating Code

Given a tiled tree:

- Postorder treewalk, with node-dependent order for children
- Emit code sequence for tiles, in order
- Tie boundaries together with register names



- Tile 6 uses registers produced by tiles 1 & 5
- Tile 6 emits “`store rtile 5 ⇒ rtile 1`”
- Can incorporate a “real” register allocator or just use virtual registers

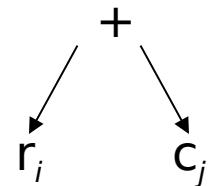
So, What's Hard About This?

Finding the matches to tile the tree

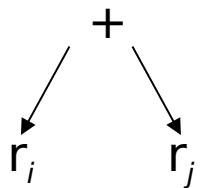
- Compiler writer connects operation trees to AST subtrees
 - Encode tree syntax, in linear form
 - Provides a set of rewrite rules
 - Associated with each is a code template

Notation

To describe these trees, we need a concise notation



$+(r_i, c_j)$



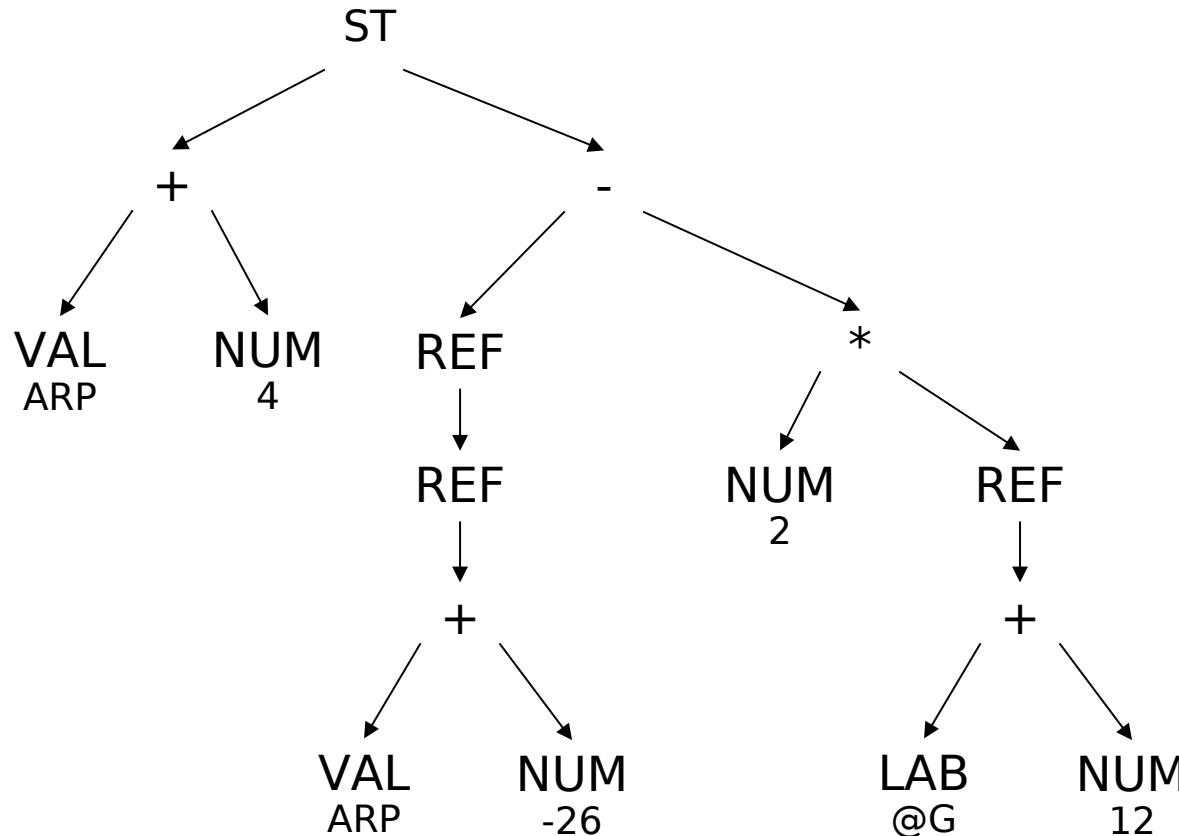
$+(r_i, r_j)$



Linear prefix form

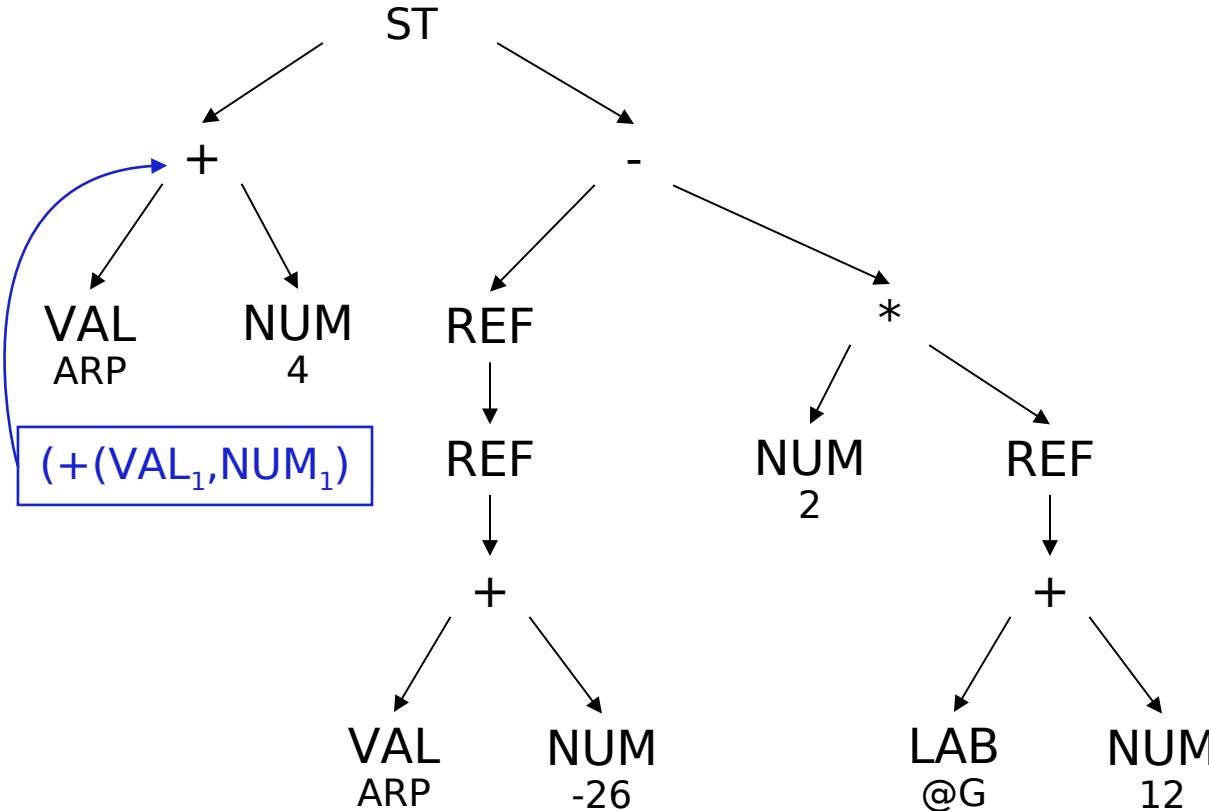
Notation

To describe these trees, we need a concise notation



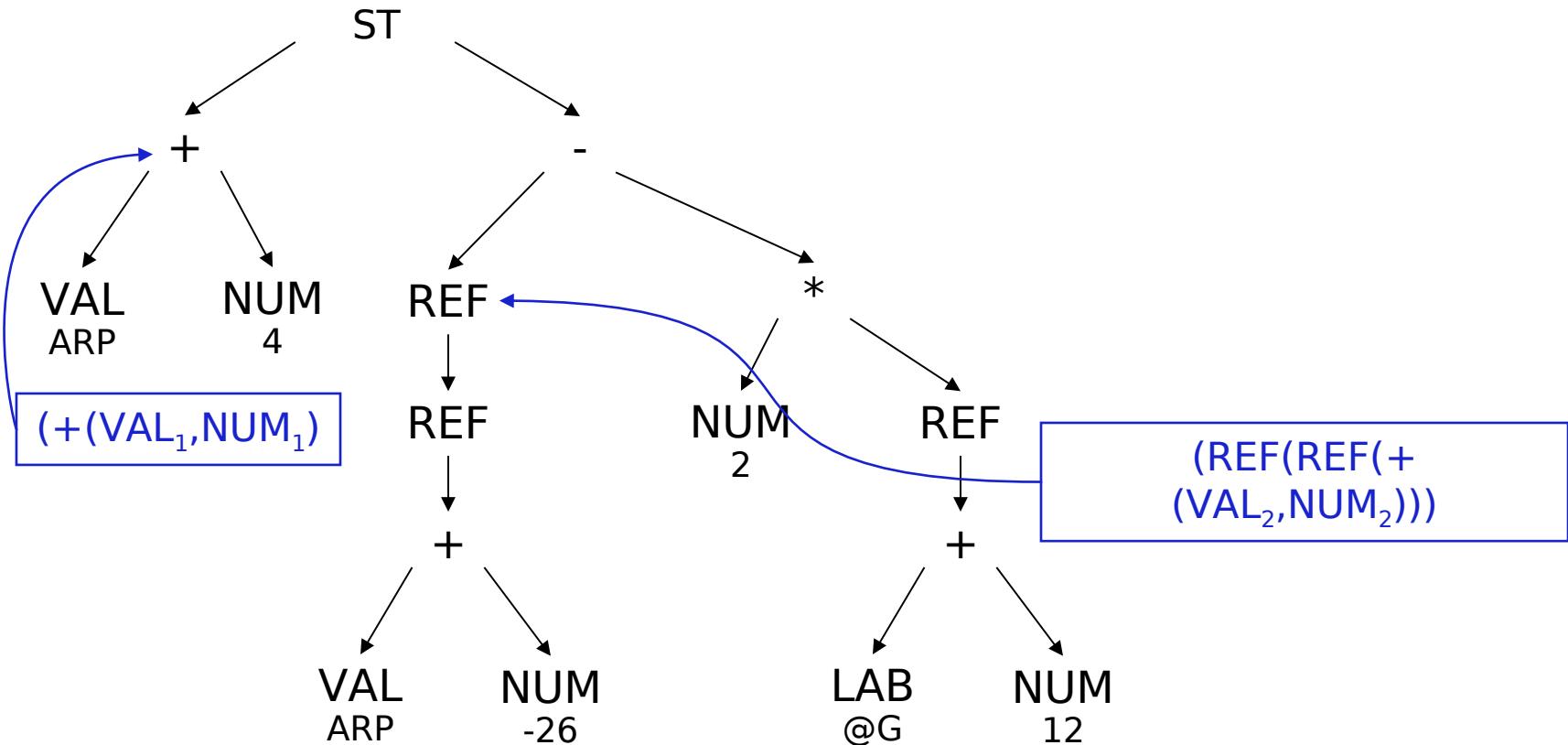
Notation

To describe these trees, we need a concise notation



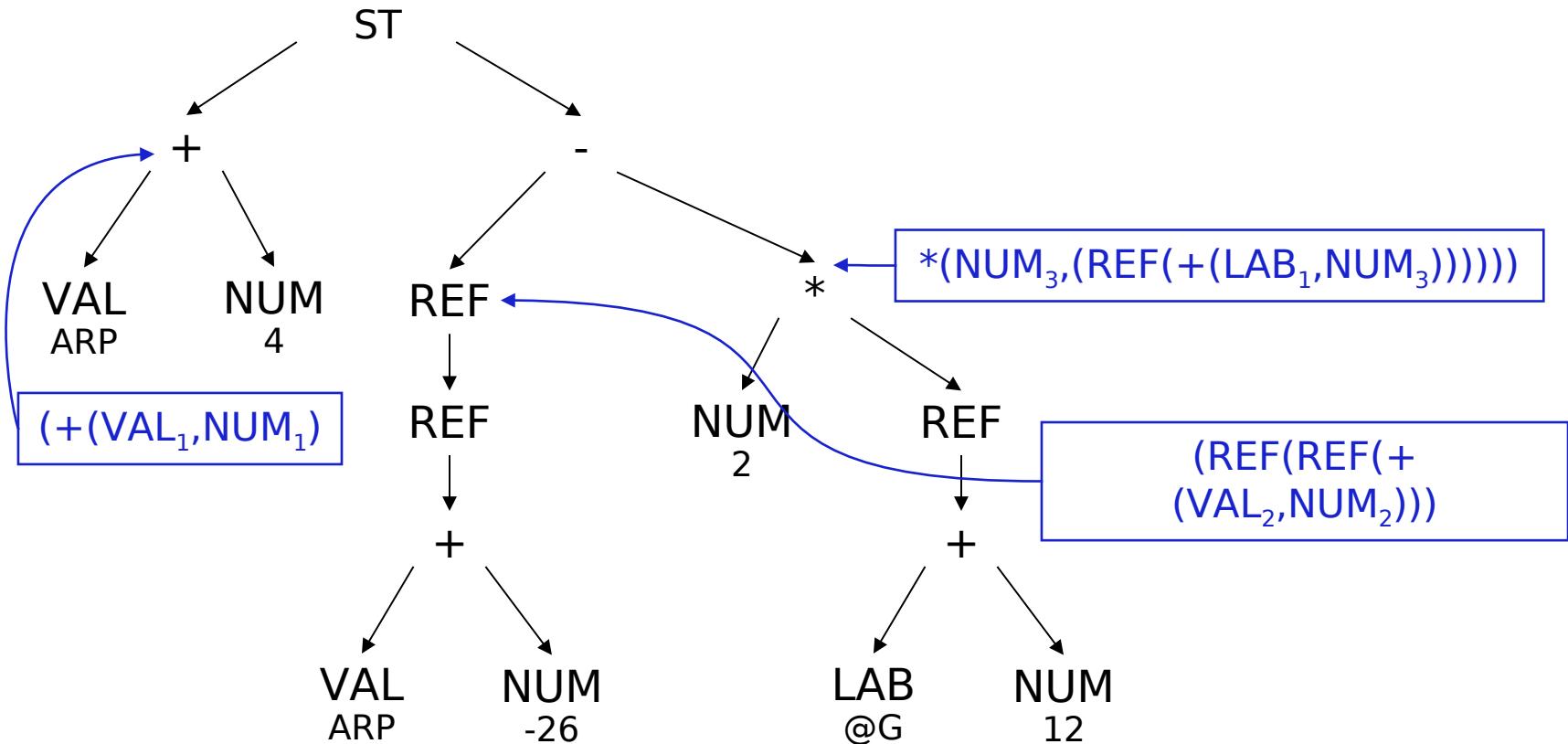
Notation

To describe these trees, we need a concise notation



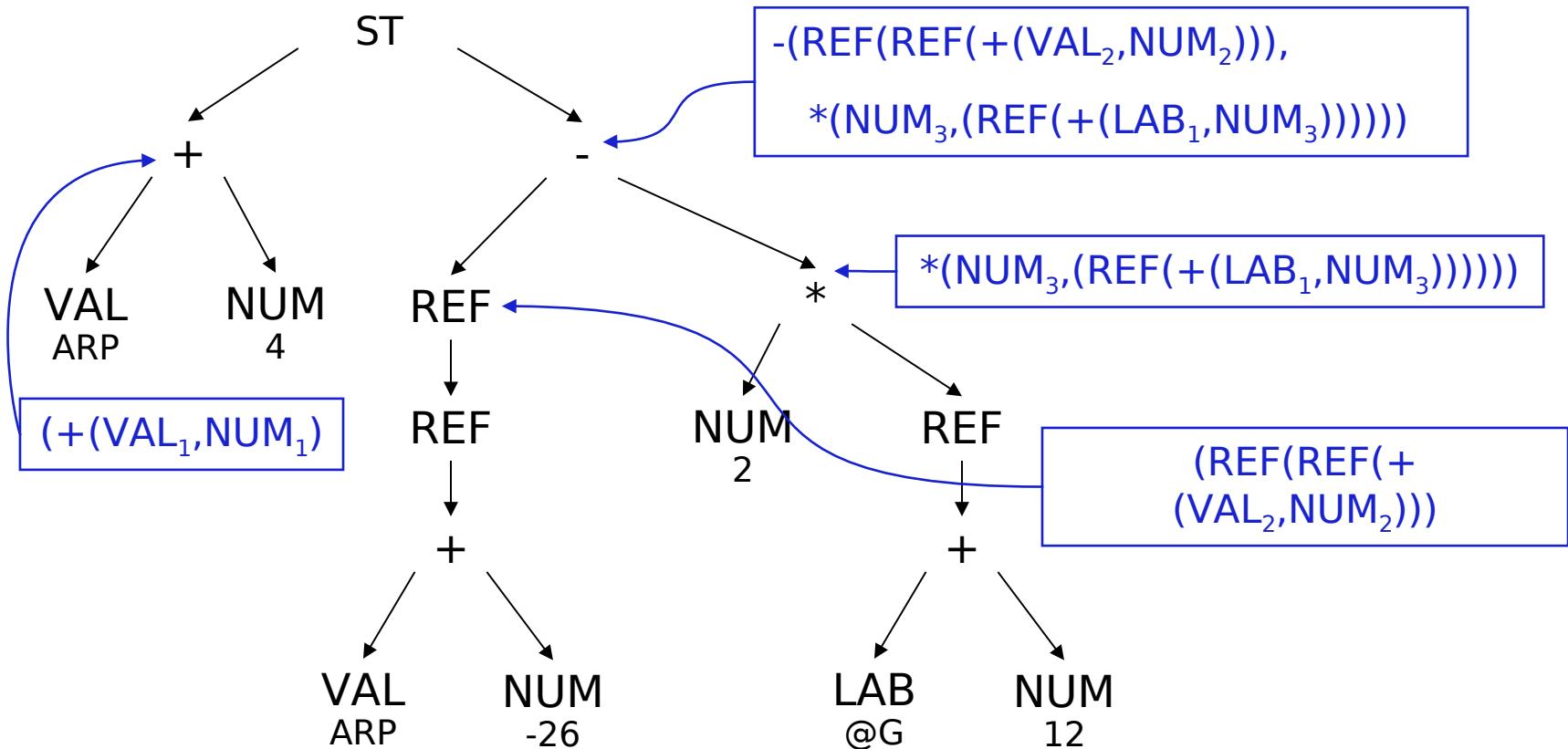
Notation

To describe these trees, we need a concise notation



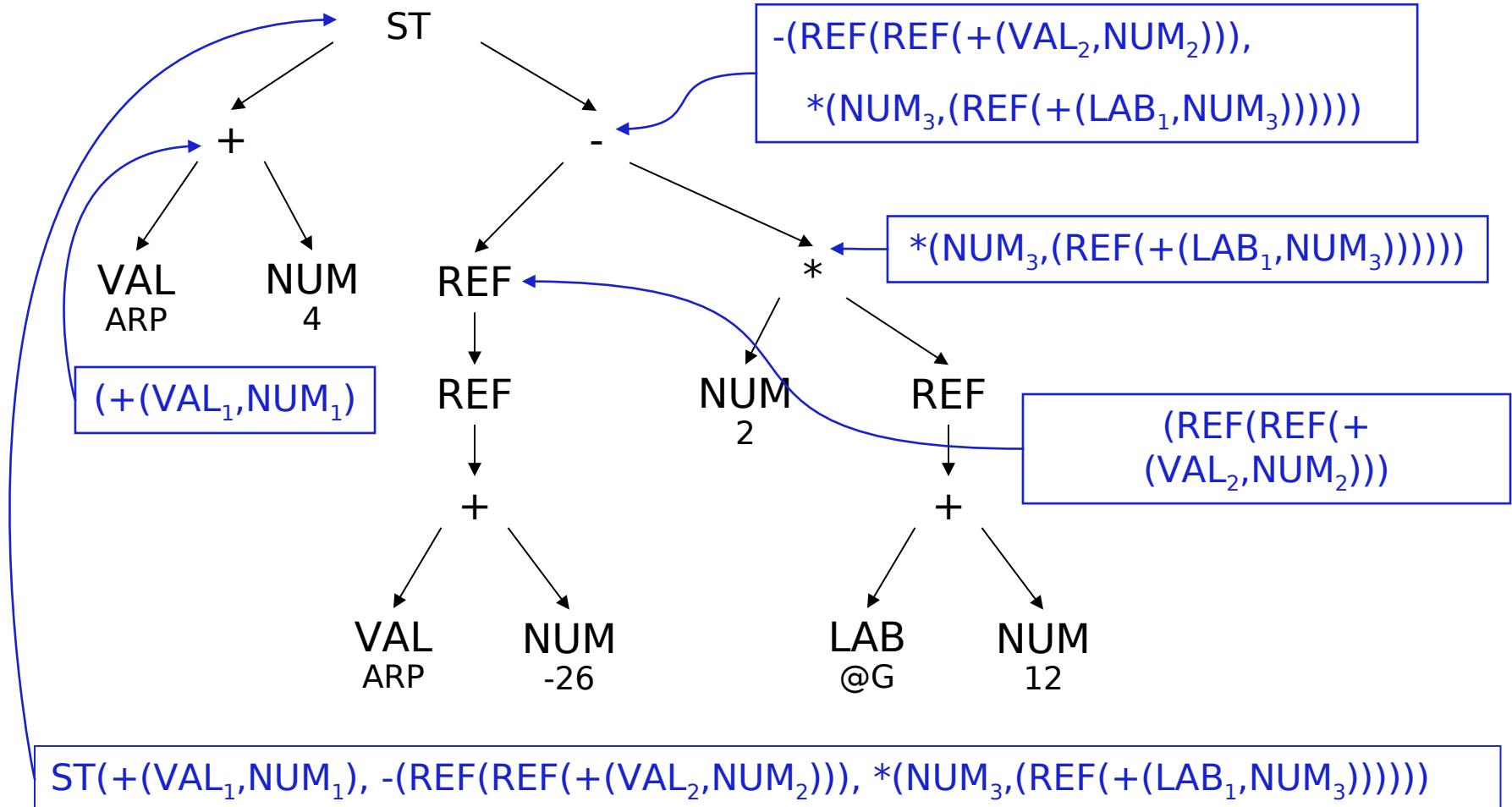
Notation

To describe these trees, we need a concise notation



Notation

To describe these trees, we need a concise notation



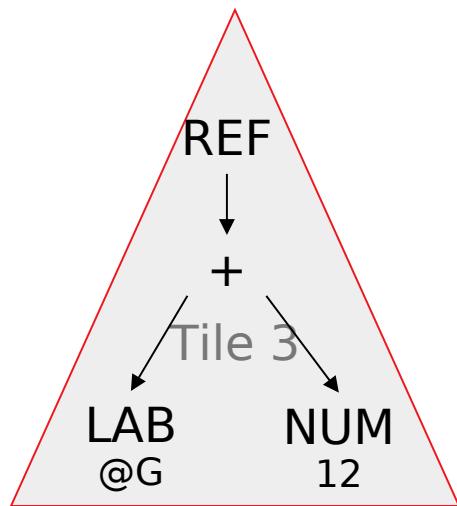
Rewrite rules: LL Integer AST into ILOC

| | Rule | Cost | Template |
|-----|--|-------------|---|
| 1 | Goal → Assign | 0 | |
| 2 | Assign → ST(Reg ₁ ,Reg ₂) | 3 | store r ₂ ⇒ r ₁ |
| 3 | Assign → ST(+ (Reg ₁ ,Reg ₂),Reg ₃) | 3 | storeAO r ₃ ⇒ r ₁ , r ₂ |
| 4 | Assign → ST(+ (Reg ₁ ,NUM ₂),Reg ₃) | 3 | storeAI r ₃ ⇒ r ₁ , n ₂ |
| 5 | Assign → ST(+ (NUM ₁ ,Reg ₂),Reg ₃) | 3 | storeAI r ₃ ⇒ r ₂ , n ₁ |
| 6 | Reg → LAB ₁ | 1 | loadI l ₁ ⇒ r _{new} |
| 7 | Reg → VAL ₁ | 0 | |
| 8 | Reg → NUM ₁ | 1 | loadI n ₁ ⇒ r _{new} |
| 9 | Reg → REF(Reg ₁) | 10 | load r ₁ ⇒ r _{new} |
| 10 | Reg → REF(+ (Reg ₁ ,Reg ₂)) | 10 | loadAO r ₁ , r ₂ ⇒ r _{new} |
| 11 | Reg → REF(+ (Reg ₁ ,NUM ₂)) | 10 | loadAI r ₁ , n ₂ ⇒ r _{new} |
| 12 | Reg → REF(+ (NUM ₁ ,Reg ₂)) | 10 | loadAI r ₂ , n ₁ ⇒ r _{new} |
| 13 | Reg → REF(+ (Reg ₁ ,Lab ₂)) | 10 | loadAI r ₁ , l ₂ ⇒ r _{new} |
| 14 | Reg → REF(+ (Lab ₁ ,Reg ₂)) | 10 | loadAI r ₂ , l ₁ ⇒ r _{new} |
| 15 | Reg → + (Reg ₁ ,Reg ₂) | 1 | addI r ₁ , r ₂ ⇒ r _{new} |
| 16 | Reg → + (Reg ₁ ,NUM ₂) | 1 | addI r ₁ , n ₂ ⇒ r _{new} |
| 17 | Reg → + (NUM ₁ ,Reg ₂) | 1 | addI r ₂ , n ₁ ⇒ r _{new} |
| 18 | Reg → + (Reg ₁ ,Lab ₂) | 1 | addI r ₁ , l ₂ ⇒ r _{new} |
| 19 | Reg → + (Lab ₁ ,Reg ₂) | 1 | addI r ₂ , l ₁ ⇒ r _{new} |
| 20 | Reg → - (NUM ₁ ,Reg ₂) | 1 | subI r ₂ , n ₁ ⇒ r _{new} |
| ... | ... | ... | ... |

A real set of rules would cover more than signed integers ...

So, What's Hard About This?

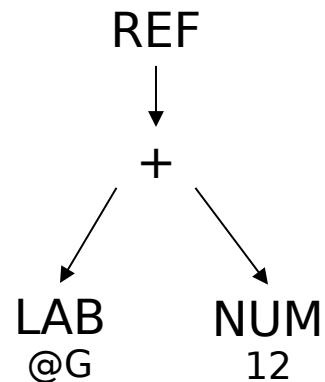
Consider tile 3 in our example



So, What's Hard About This?

Consider tile 3 in our example

What rules match tile 3?

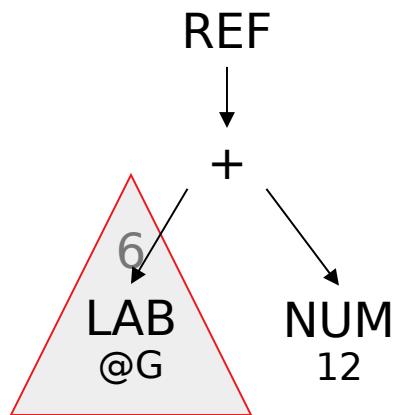


So, What's Hard About This?

Consider tile 3 in our example

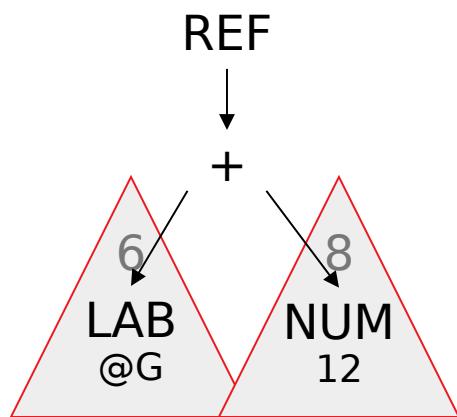
What rules match tile 3?

6: Reg → LAB₁ tiles the lower left node



So, What's Hard About This?

Consider tile 3 in our example



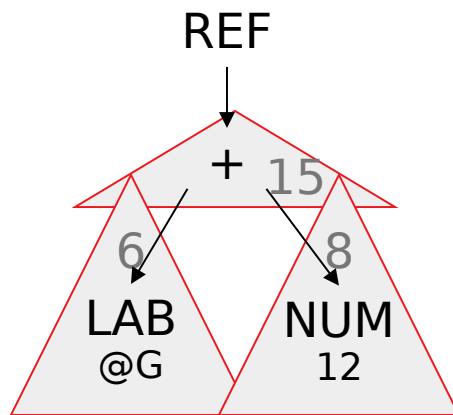
What rules match tile 3?

6: Reg \rightarrow LAB₁ tiles the lower left node

8: Reg \rightarrow NUM₁ tiles the bottom right node

So, What's Hard About This?

Consider tile 3 in our example



What rules match tile 3?

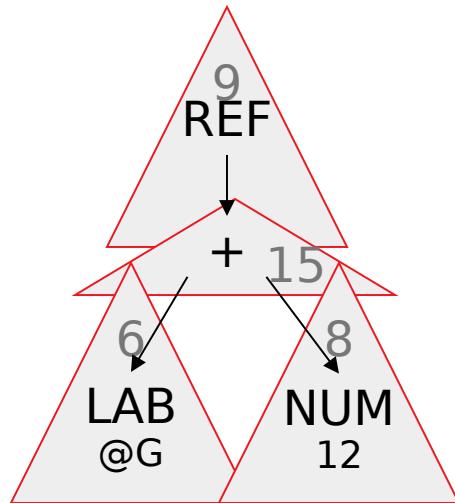
6: Reg → LAB₁ tiles the lower left node

8: Reg → NUM₁ tiles the bottom right node

15: Reg → + (Reg₁,Reg₂) tiles the + node

So, What's Hard About This?

Consider tile 3 in our example

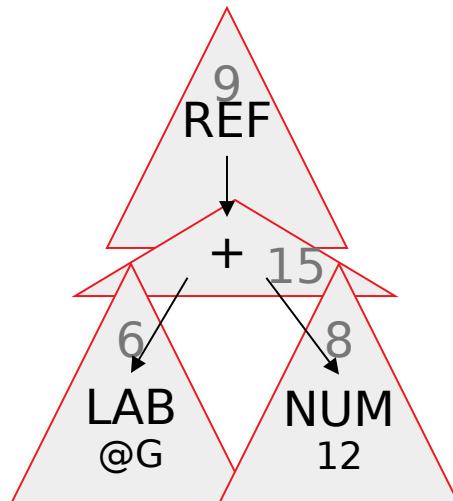


What rules match tile 3?

- 6: $\text{Reg} \rightarrow \text{LAB}_1$ tiles the lower left node
- 8: $\text{Reg} \rightarrow \text{NUM}_1$ tiles the bottom right node
- 15: $\text{Reg} \rightarrow + (\text{Reg}_1, \text{Reg}_2)$ tiles the + node
- 9: $\text{Reg} \rightarrow \text{REF}(\text{Reg}_1)$ tiles the REF

So, What's Hard About This?

Consider tile 3 in our example



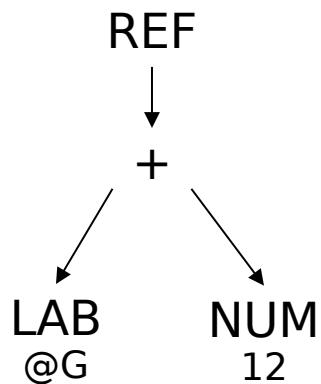
What rules match tile 3?

- 6: Reg \rightarrow LAB₁ tiles the lower left node
- 8: Reg \rightarrow NUM₁ tiles the bottom right node
- 15: Reg \rightarrow + (Reg₁,Reg₂) tiles the + node
- 9: Reg \rightarrow REF(Reg₁) tiles the REF

We denote this match as $<6,8,15,9>$
Of course, it implies $<8,6,15,9>$
Both have a cost of 13

Finding matches

Many Sequences Match Our Subtree



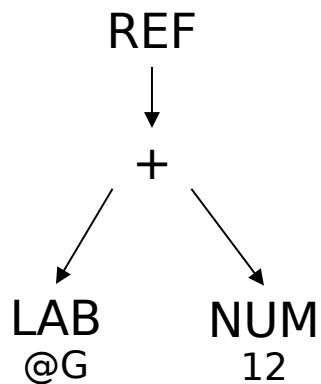
| Cost | Sequences | | | |
|------|-----------|----------|--------|--------|
| 11 | 6,11 | 8,14 | | |
| 12 | 6,8,10 | 8,6,10 | 6,16,9 | 8,19,9 |
| 13 | 6,8,15,9 | 8,6,15,9 | | |

In general, we want the low cost sequence

- Each unit of cost is an operation (1 cycle)
- We should favour short sequences

Finding matches

Low Cost Matches



| Sequences with Cost of 11 | |
|---|---|
| 6: Reg → LAB ₁ | loadI @G ⇒ r _i |
| 11: Reg → REF(+ (Reg ₁ , NUM ₂)) | loadAI r _i , 12 ⇒ r _j |
| 8: Reg → NUM ₁ | loadI 12 ⇒ r _i |
| 14: Reg → REF(+ (LAB ₁ , Reg ₂)) | loadAI r _i , @G ⇒ r _j |

These two are equivalent in cost

6,11 might be better, because @G may be longer than the immediate field

Tiling the Tree

- Assume each rule implements **one** operator
- Assume operator takes 0, 1, or 2 operands
- Need to “massage” the rules, e.g.:

| | Rule | \$ | Template |
|-----|---|-----|--|
| ... | ... | ... | ... |
| 10 | $\text{Reg} \rightarrow \text{REF}(+ (\text{Reg}_1, \text{Reg}_2))$ | 10 | $\text{loadAO } r_1, r_2 \Rightarrow r_{\text{new}}$ |
| ... | ... | ... | ... |

becomes:

| | Rule | \$ | Template |
|-----|--|-----|--|
| ... | ... | ... | ... |
| 10 | $\text{Reg} \rightarrow \text{REF(R10a)}$ | 10 | $\text{loadAO } r_1, r_2 \Rightarrow r_{\text{new}}$ |
| 10a | $\text{R10a} \rightarrow + (\text{Reg}_1, \text{Reg}_2)$ | 0 | |
| ... | ... | ... | ... |

Algorithm to tile the tree

```
Tile(n)
  Label(n) ← Ø

  if n is a leaf
    Label(n) = {n → n} // dummy rule
    Label(n) += all rules where rhs == n

  else if n is a unary node
    Tile(child(n))
    foreach rule r that matches n's op
      if (child(r) in Label(child(n)))
        Label(n) += r

  else // n is binary node
    Tile(left(n))
    Tile(right(n))
    foreach rule r that matches n's op
      if (left(rhs(r)) in lhs(Label(left(n))) &
           right(rhs(r)) in lhs(Label(right(n)))))
        Label(n) += r
```

Notes:

- child, left and right refer to the children of an AST node or of the rhs of a rule

Tiling the Tree

```

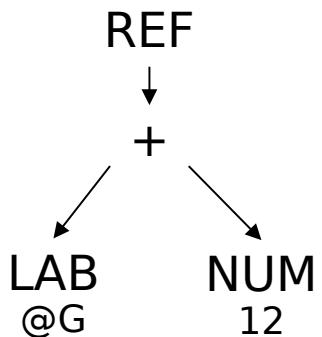
Tile(n)
Label(n) ← Ø

if n is a leaf
    Label(n) = {n → n} // dummy rule
    Label(n) += all rules where rhs == n

else if n is a unary node
    Tile(child(n))
    foreach rule r that matches n's op
        if (child(r) in Label(child(n)))
            Label(n) += r

else // n is binary node
    Tile(left(n))
    Tile(right(n))
    foreach rule r that matches n's op
        if (left(rhs(r)) in lhs(Label(left(n))) &
            right(rhs(r)) in lhs(Label(right(n)))))
            Label(n) += r

```



| | Rule | \$ | Template |
|-----|---|-----|---|
| ... | ... | ... | ... |
| 6 | Reg → LAB ₁ | 1 | loadI l ₁ ⇒ r _{new} |
| 7 | Reg → VAL ₁ | 0 | |
| 8 | Reg → NUM ₁ | 1 | loadI n ₁ ⇒ r _{new} |
| 9 | Reg → REF(Reg ₁) | 10 | load r ₁ ⇒ r _{new} |
| 10 | Reg → REF(R10a) | 10 | loadAO r ₁ , r ₂ ⇒ r _{new} |
| 10a | R10a → + (Reg ₁ , Reg ₂) | 0 | |
| 11 | Reg → REF(R11a) | 10 | loadAI r ₁ , n ₂ ⇒ r _{new} |
| 11a | R11a → + (Reg ₁ , NUM ₂) | 0 | |
| 12 | Reg → REF(R12a) | 10 | loadAI r ₂ , n ₁ ⊖ r _{new} |
| 12a | R12a → + (NUM ₁ , Reg ₂) | 0 | |
| 13 | Reg → REF(R13a) | 10 | loadAI r ₁ , l ₂ ⊖ r _{new} |
| 13a | R13a → + (Reg ₁ , Lab ₂) | 0 | |
| 14 | Reg → REF(R14a) | 10 | loadAI r ₂ , l ₁ ⊖ r _{new} |
| 14a | R14a → + (Lab ₁ , Reg ₂) | 0 | |
| 15 | Reg → + (Reg ₁ , Reg ₂) | 1 | addI r ₁ , r ₂ ⇒ r _{new} |
| 16 | Reg → + (Reg ₁ , NUM ₂) | 1 | addI r ₁ , n ₂ ⇒ r _{new} |
| 17 | Reg → + (NUM ₁ , Reg ₂) | 1 | addI r ₂ , n ₁ ⇒ r _{new} |
| 18 | Reg → + (Reg ₁ , Lab ₂) | 1 | addI r ₁ , l ₂ ⇒ r _{new} |
| 19 | Reg → + (Lab ₁ , Reg ₂) | 1 | addI r ₂ , l ₁ ⇒ r _{new} |
| ... | ... | ... | ... |

Label(Ref) =

Label(+) =

Label(Lab) =

Label(Num) =

Tiling the Tree

```
Tile(n)
Label(n) ← Ø

if n is a leaf
    Label(n) = {n → n} // dummy rule
    Label(n) += all rules where rhs == n

else if n is a unary node
    Tile(child(n))
    foreach rule r that matches n's op
        if (child(r) in Label(child(n)))
            Label(n) += r

else // n is binary node
    Tile(left(n))
    Tile(right(n))
    foreach rule r that matches n's op
        if (left(rhs(r)) in lhs(Label(left(n))) &
            right(rhs(r)) in lhs(Label(right(n)))))
            Label(n) += r
```

This algorithm

- Finds all matches in rule set
- Labels node n with that set
- Can keep lowest cost match at each point for each type of nodes
→ Dynamic programming
- Spends its time in the two matching loops

The Big Picture

- Tree patterns represent AST and ASM
- Can use matching algorithms to find low-cost tiling of AST
- Can turn a tiling into code using templates for matched rules
- Techniques (& tools) exist to do this efficiently

| | |
|--|--|
| Hand-coded matcher | Lots of work |
| Encode matching as an automaton | $O(1)$ cost per node Tools like BURS (bottom-up rewriting system), BURG |
| Use parsing techniques | Uses known technology Very ambiguous grammars |
| Linearize tree into string and use string searching algorithm (Aho-Corasick) | Finds all matches |

Next Lecture

- Instruction scheduling