# Instruction Selection:
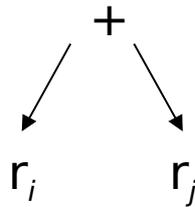# Tree-pattern matching

EaC-11.3

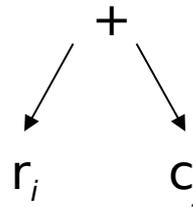# The Concept

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 ILOC add operators
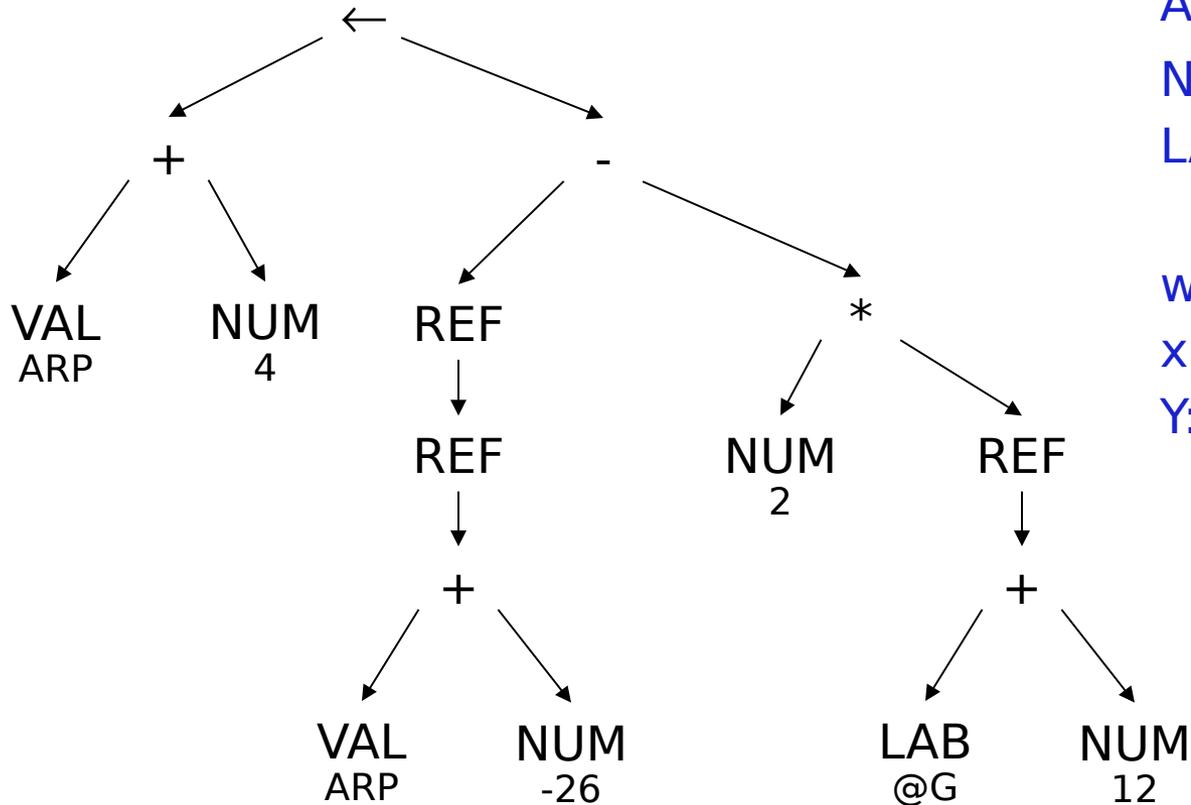


$$\text{add } r_i, r_j \Rightarrow r_k \qquad\qquad \text{addI } r_i, c_j \Rightarrow r_k$$

Operation trees

If we can match these "pattern trees" against IR trees, …

# The Concept

Low-level AST for w ← x - 2 * y



ARP:   $r_{arp}$
NUM: constant
LAB:   ASM label

w: at ARP+4
x: at ARP-26
Y: at @G+12

# The Concept

Low-level AST for w ← x - 2 * y



ARP:   $r_{arp}$
NUM: constant
LAB:   ASM label

w: at ARP+4
x: at ARP-26
Y: at @G+12

# 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



GETS

$-(REF(REF(+(VAL_2,NUM_2))),$
$\quad *(NUM_3,(REF(+(LAB_1,NUM_3))))))$

+

-

$*(NUM_3,(REF(+(LAB_1,NUM_3))))))$

VAL
ARP

NUM
4

REF

*

$(+(VAL_1,NUM_1)$

REF

NUM
2

REF

$(REF(REF(+(VAL_2,NUM_2)))$

+

+

VAL
ARP

NUM
-26

LAB
@G

NUM
12

$GETS(+(VAL_1,NUM_1), -(REF(REF(+(VAL_2,NUM_2))), *(NUM_3,(REF(+(LAB_1,NUM_3))))))$

# Tree-pattern matching

Goal is to "tile" AST with operation trees

- A tiling is collection of <*ast,op* > pairs
  - → *ast* is a node in the AST
  - → *op* is an operation tree
  - → <*ast*, *op* > 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
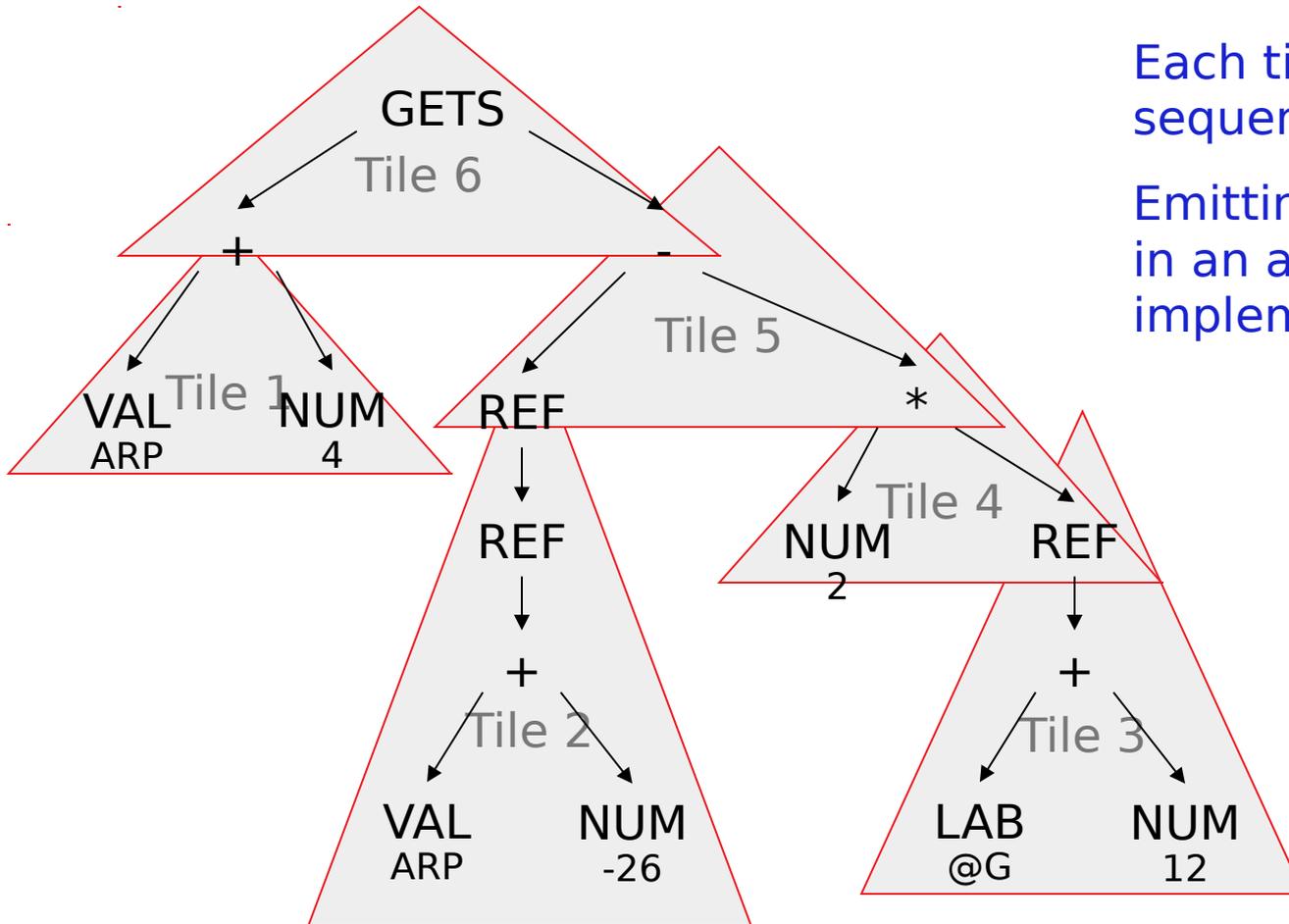  - → <*ast, op*> ∈ tiling means *ast* is also covered by a leaf in another operation tree in the tiling, unless it is the root
  - → Where two operation trees meet, they must be compatible (expect the value in the same location)

# Tiling the Tree



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
  - → Right child of GETS before its left child
  - → Might impose "most demanding first" rule ...            (*Sethi* )

- 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` $r_{tile\ 5}$ ⇒ $r_{tile\ 1}$"
  - → Can incorporate a "real" allocator or can use "NextRegister++"

# So, What's Hard About This?

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

# Rewrite rules: LL Integer AST into ILOC

| | Rule | Cost | Template |
|---|---|---|---|
| 1 | Goal $\rightarrow$ Assign | 0 | |
| 2 | Assign $\rightarrow$ GETS(Reg$_1$,Reg$_2$) | 1 | `store`      r$_2$ $\Rightarrow$ r$_1$ |
| 3 | Assign $\rightarrow$ GETS(+(Reg$_1$,Reg$_2$),Reg$_3$) | 1 | `storeAO`  r$_3$ $\Rightarrow$ r$_1$,r$_2$ |
| 4 | Assign $\rightarrow$ GETS(+(Reg$_1$,NUM$_2$),Reg$_3$) | 1 | `storeAI`  r$_3$ $\Rightarrow$ r$_1$,n$_2$ |
| 5 | Assign $\rightarrow$ GETS(+(NUM$_1$,Reg$_2$),Reg$_3$) | 1 | `storeAI`  r$_3$ $\Rightarrow$ r$_2$,n$_1$ |
| 6 | Reg $\rightarrow$ LAB$_1$ | 1 | `loadI`      l$_1$ $\Rightarrow$ r$_{new}$ |
| 7 | Reg $\rightarrow$ VAL$_1$ | 0 | |
| 8 | Reg $\rightarrow$ NUM$_1$ | 1 | `loadI`      n$_1$ $\Rightarrow$ r$_{new}$ |
| 9 | Reg $\rightarrow$ REF(Reg$_1$) | 1 | `load`      r$_1$ $\Rightarrow$ r$_{new}$ |
| 10 | Reg $\rightarrow$ REF(+ (Reg$_1$,Reg$_2$)) | 1 | `loadAO` r$_1$,r$_2$ $\Rightarrow$ r$_{new}$ |
| 11 | Reg $\rightarrow$ REF(+ (Reg$_1$,NUM$_2$)) | 1 | `loadAI` r$_1$,n$_2$ $\Rightarrow$ r$_{new}$ |
| 12 | Reg $\rightarrow$ REF(+ (NUM$_1$,Reg$_2$)) | 1 | `loadAI` r$_2$,n$_1$ $\Rightarrow$ r$_{new}$ |

# Rewrite rules: LL Integer AST into ILOC (*part II*)

| | Rule | Cost | Template |
|---|---|---|---|
| 13 | Reg $\rightarrow$ + (Reg$_1$,Reg$_2$) | 1 | `add` $\quad$ r$_1$,r$_2$ $\Rightarrow$ r$_{new}$ |
| 14 | Reg $\rightarrow$ + (Reg$_1$,NUM$_2$) | 1 | `addI` $\quad$ r$_1$,n$_2$ $\Rightarrow$ r$_{new}$ |
| 15 | Reg $\rightarrow$ + (NUM$_1$,Reg$_2$) | 1 | `addI` $\quad$ r$_2$,n$_1$ $\Rightarrow$ r$_{new}$ |
| 16 | Reg $\rightarrow$ - (Reg$_1$,Reg$_2$) | 1 | `sub` $\quad$ r$_1$,r$_2$ $\Rightarrow$ r$_{new}$ |
| 17 | Reg $\rightarrow$ - (Reg$_1$,NUM$_2$) | 1 | `subI` $\quad$ r$_1$,n$_2$ $\Rightarrow$ r$_{new}$ |
| 18 | Reg $\rightarrow$ - (NUM$_1$,Reg$_2$) | 1 | `rsubI` $\quad$ r$_2$,n$_1$ $\Rightarrow$ r$_{new}$ |
| 19 | Reg $\rightarrow$ x (Reg$_1$,Reg$_2$) | 1 | `mult` $\quad$ r$_1$,r$_2$ $\Rightarrow$ r$_{new}$ |
| 20 | Reg $\rightarrow$ x (Reg$_1$,NUM$_2$) | 1 | `multI` $\quad$ r$_1$,n$_2$ $\Rightarrow$ r$_{new}$ |
| 21 | Reg $\rightarrow$ x (NUM$_1$,Reg$_2$) | 1 | `multI` $\quad$ r$_2$,n$_1$ $\Rightarrow$ r$_{new}$ |

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

# So, What's Hard About This?

Need an algorithm to AST subtrees with the rules

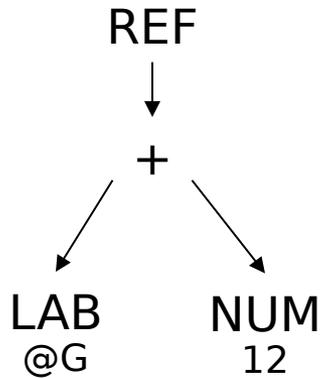Consider tile 3 in our example

REF

↓

+

Tile 3

LAB          NUM
@G           12

# So, What's Hard About This?

Need an algorithm to AST subtrees with the rules

Consider tile 3 in our example

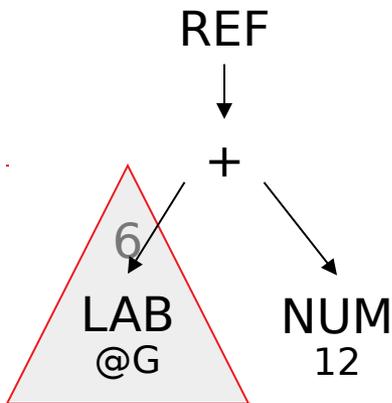What rules match tile 3?

REF

↓

+

LAB     NUM
@G       12

# So, What's Hard About This?

Need an algorithm to AST subtrees with the rules

Consider tile 3 in our example

What rules match tile 3?

6:  Reg $\rightarrow$ LAB$_1$ tiles the lower left node

REF

$\downarrow$

+

6

LAB
@G

NUM
12

# So, What's Hard About This?

Need an algorithm to AST subtrees with the rules

Consider tile 3 in our example

REF

$\downarrow$

+

6     8

LAB     NUM
@G      12

What rules match tile 3?

6:  Reg $\rightarrow$ LAB$_1$ tiles the lower left node

8: Reg $\rightarrow$ NUM$_1$ tiles the bottom right node

# So, What's Hard About This?

Need an algorithm to AST subtrees with the rules
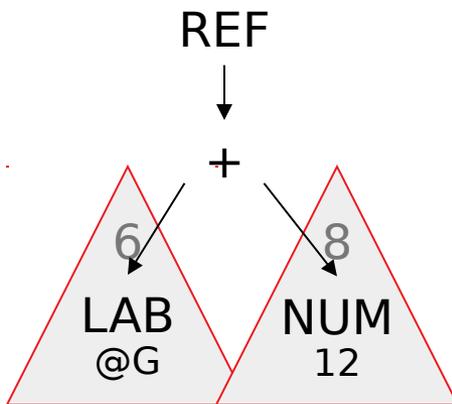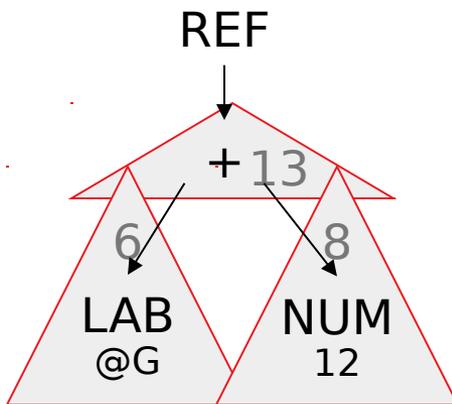
Consider tile 3 in our example

What rules match tile 3?

6:  Reg $\rightarrow$ LAB$_1$ tiles the lower left node

8: Reg $\rightarrow$ NUM$_1$ tiles the bottom right node

13: Reg $\rightarrow$ + (Reg$_1$,Reg$_2$) tiles the + node

REF

+ 13

6

8

LAB
@G

NUM
12

# So, What's Hard About This?

Need an algorithm to AST subtrees with the rules

Consider tile 3 in our example



What rules match tile 3?

6: Reg $\rightarrow$ LAB$_1$ tiles the lower left node

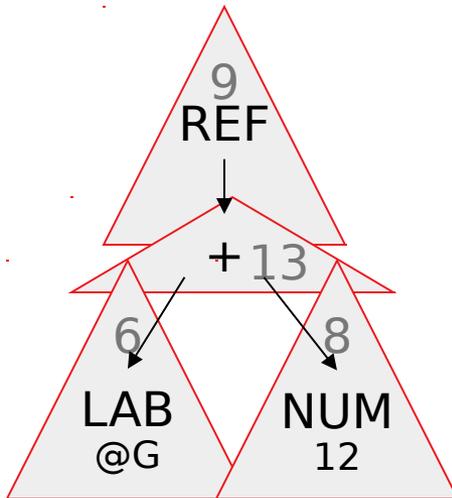8: Reg $\rightarrow$ NUM$_1$ tiles the bottom right node

13: Reg $\rightarrow$ + (Reg$_1$,Reg$_2$) tiles the + node

9: Reg $\rightarrow$ REF(Reg$_1$) tiles the REF

# So, What's Hard About This?

Need an algorithm to AST subtrees with the rules

Consider tile 3 in our example

What rules match tile 3?

6:  Reg $\rightarrow$ LAB$_1$ tiles the lower left node
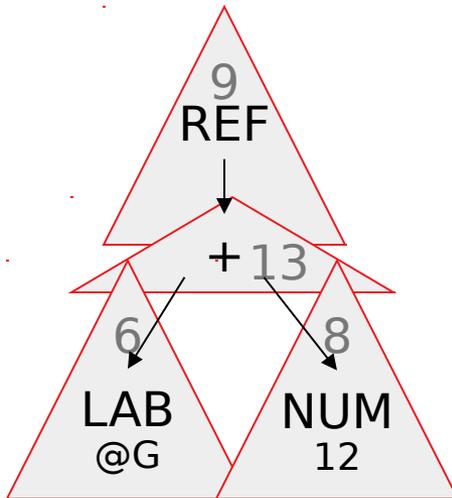
8: Reg $\rightarrow$ NUM$_1$ tiles the bottom right node

13: Reg $\rightarrow$ + (Reg$_1$,Reg$_2$) tiles the + node

9:  Reg $\rightarrow$ REF(Reg$_1$) tiles the REF

We denote this match as <6,8,13,9>
Of course, it implies <8,6,13,9>
Both have a cost of 4

# Finding matches

Many Sequences Match Our Subtree

| Cost | Sequences | | | |
|------|-----------|-----------|--------|--------|
| 2 | 6,11 | 8,12 | | |
| 3 | 6,8,10 | 8,6,10 | 6,14,9 | 8,15,9 |
| 4 | 6,8,13,9 | 8,6,13,9 | | |

REF

↓

+

↙ ↘

LAB     NUM
@G      12

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

REF
$\downarrow$
+
LAB    NUM
@G      12

| Sequences with Cost of 2 | |
|---|---|
| 6: Reg $\rightarrow$ LAB$_1$<br>11: Reg $\rightarrow$ REF(+(Reg$_1$,NUM$_2$)) | `loadI  @G    ` $\Rightarrow$ `r`$_i$<br>`loadAI r`$_i$`,12` $\Rightarrow$ `r`$_j$ |
| 8: Reg $\rightarrow$ NUM$_1$<br>12: Reg $\rightarrow$ REF(+(NUM$_1$,Reg$_2$)) | `loadI 12    ` $\Rightarrow$ `r`$_i$<br>`loadAI r`$_i$`,@G` $\Rightarrow$ `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

Still need an algorithm

- Assume each rule implements one operator
- Assume operator takes 0, 1, or 2 operands

Now, …

# Tiling the Tree

*Tile(n)*
  *Label(n) ← Ø*
  *if n has two children then*
    *Tile (left child of n)*
    *Tile (right child of n)*
    *for each rule r that implements n*
      *if (left(r) ∈ Label(left(n)) and*
        *(right(r) ∈ Label(right(n))*
      *then Label(n) ← Label(n) ∪ { r }*

<span style="color:blue">Match binary nodes against binary rules</span>

  *else if n has one child*
    *Tile(child of n)*
    *for each rule r that implements n*
      *if (left(r) ∈ Label(child(n))*
        *then Label(n) ← Label(n) ∪ { r }*

<span style="color:blue">Match unary nodes against unary rules</span>

  *else  /\* n is a leaf \*/*
    *Label(n) ← {all  rules that implement n }*

<span style="color:blue">Handle leaves with lookup in rule table</span>

# Tiling the Tree

*Tile(n)*
  *Label(n) ← Ø*
  *if n has two children then*
      *Tile (left child of n)*
      *Tile (right child of n)*
      *for each rule r that implements n*
          *if (left(r) ∈ Label(left(n)) and*
              *(right(r) ∈ Label(right(n))*
              *then Label(n) ← Label(n) ∪ { r }*

  *else if n has one child*
      *Tile(child of n)*
      *for each rule r that implements n*
          *if (left(r) ∈ Label(child(n))*
              *then Label(n) ← Label(n) ∪ { r }*

  *else  /* n is a leaf */*
      *Label(n) ← {all  rules that implement n }*

This algorithm
- Finds all matches in rule set
- Labels node n with that set
- Can keep lowest cost match at each point
- Leads to a notion of local optimality — lowest cost at each point
- Spends its time in the two matching loops

# Tiling the Tree

*Tile(n)*
   *Label(n) ← Ø*
  *if n has two children then*
     *Tile (left child of n)*
     *Tile (right child of n)*
     *for each rule r that implements n*
       *if (left(r) ∈ Label(left(n)) and*
        *(right(r) ∈ Label(right(n))*
       *then Label(n) ← Label(n) ∪ { r }*

  *else if n has one child*
     *Tile(child of n)*
     *for each rule r that implements n*
       *if (left(r) ∈ Label(child(n))*
       *then Label(n) ← Label(n) ∪ { r }*

  *else  /* n is a leaf */*
     *Label(n) ← {all  rules that implement n }*

Oversimplifications
1. Only handles 1 storage class
2. Must track low cost sequence in each class
3. Must choose lowest cost for subtree, across all classes

The extensions to handle these complications are pretty straightforward.

# Tiling the Tree

*Tile(n)*
  *Label(n) ← ∅*
  *if n has two children then*
      *Tile (left child of n)*
      *Tile (right child of n)*
      *for each rule r that implements n*
          *if (left(r) ∈ Label(left(n)) and*
              *(right(r) ∈ Label(right(n))*
              *then Label(n) ← Label(n) ∪ { r }*
  *else if n has one child*
      *Tile(child of n)*
      *for each rule r that implements n*
          *if (left(r) ∈ Label(child(n))*
              *then Label(n) ← Label(n) ∪ { r }*
  *else  /* n is a leaf */*
      *Label(n) ← {all  rules that implement n }*

Can turn matching code (inner loop) into a table lookup

Table can get huge and sparse

|op trees| x |labels| x |labels|
    200      x   1000  x   1000
leads to 200,000,000 entries

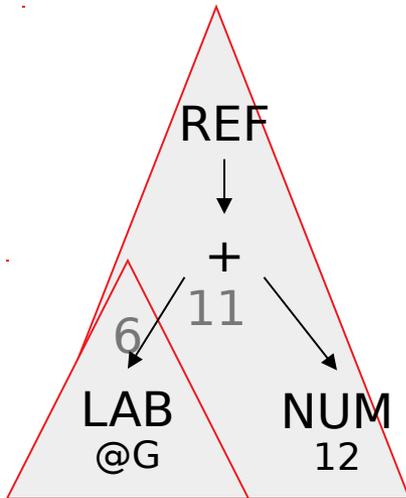Fortunately, they are quite sparse & have reasonable encodings (e.g., Chase's work)

# 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

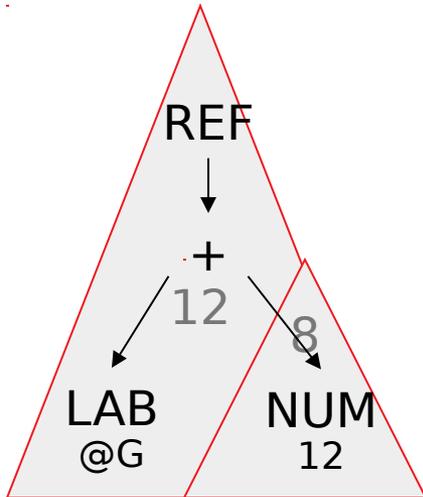| | |
|---|---|
| Hand-coded matcher like *Tile* | Avoids large sparse table<br>Lots of work |
| Encode matching as an automaton | O(1) cost per node<br>Tools like BURS (bottom-up rewriting system), BURG |
| Use parsing techniques | Uses known technology<br>Very ambiguous grammars |
| Linearize tree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

# Extra Slides Start Here

# Other Sequences



6,11

  6: $Reg \rightarrow LAB_1$

  11: $Reg \rightarrow REF( + (Reg_1, NUM_2))$

Two operator rule

# Other Sequences



REF

$\downarrow$

+

12   8

LAB
@G

NUM
12

8,12

   8:  Reg $\rightarrow$ NUM$_1$

   12: Reg $\rightarrow$ REF( + (NUM$_1$,Reg$_2$))

Two operator rule

# Other Sequences

REF

+

10

6        8

LAB        NUM
@G          12
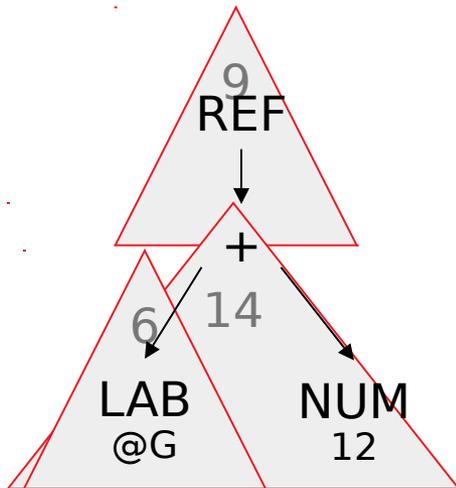
6,8,10

  6:  Reg $\rightarrow$ LAB$_1$

  8:  Reg $\rightarrow$ NUM$_1$

11: Reg $\rightarrow$ REF( + (Reg$_1$,Reg$_2$))

Two operator rule

8,6,10 looks the same

# Other Sequences



6,14,9

  6:  Reg $\to$ LAB$_1$

  14:  Reg $\to$ + (Reg$_1$,NUM$_2$)

  9: Reg $\to$ REF(Reg$_1$)

All single operator rules
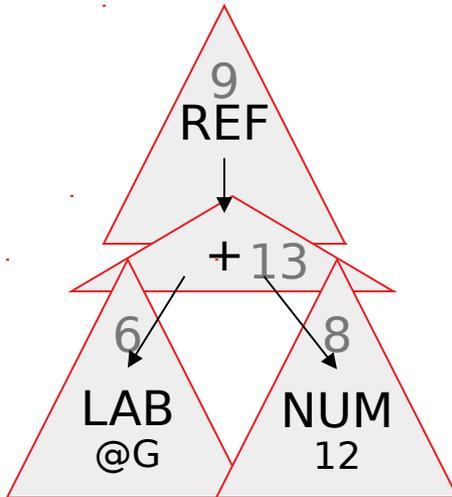
# Other Sequences



8,15,9

  8: $\text{Reg} \rightarrow \text{NUM}_1$

  15: $\text{Reg} \rightarrow + (\text{NUM}_1, \text{Reg}_2)$

  9: $\text{Reg} \rightarrow \text{REF}(\text{Reg}_1)$

All single operator rules

# Other Sequences



6,8,13,9

  6:  Reg $\rightarrow$ LAB$_1$

  8: Reg $\rightarrow$ NUM$_1$

  13: Reg $\rightarrow$ + (Reg$_1$,Reg$_2$)

  9:  Reg $\rightarrow$ REF(Reg$_1$)

All single operator rules

8,6,13,9 looks the same