



Instruction Selection, II

Tree-pattern matching

COMP 412
Fall 2005

Copyright 2005, Keith D. Cooper, Ken Kennedy & Linda Torczon, all rights reserved.
Students enrolled in Comp 412 at Rice University have explicit permission to make copies of these materials for their personal use.



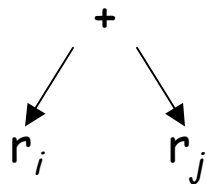
The Concept

Many compilers use tree-structured IRs

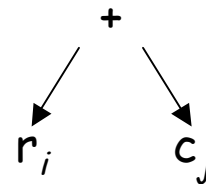
- Abstract syntax trees generated in the parser
- Trees or DAGs for expressions

These systems might well use trees to represent target ISA

Consider the ILOC add operations



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



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

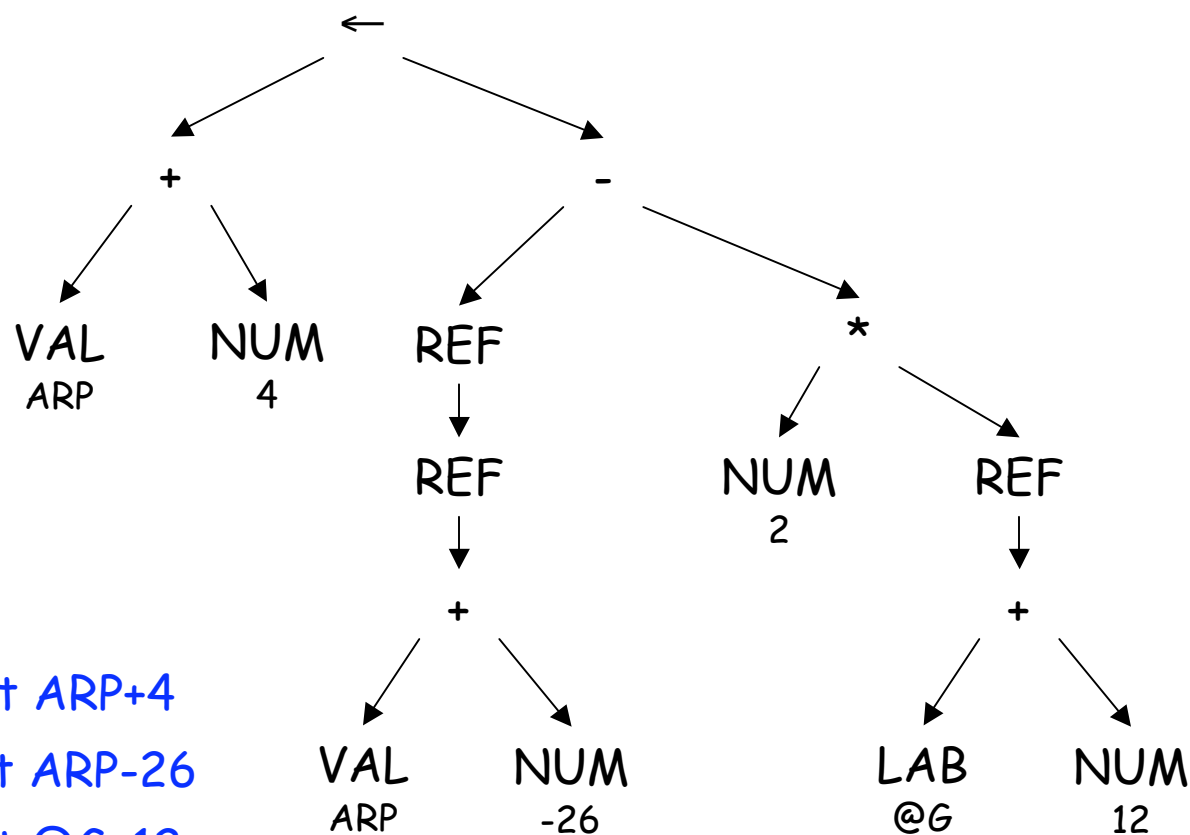
Operation trees
or pattern trees

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



The Concept

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

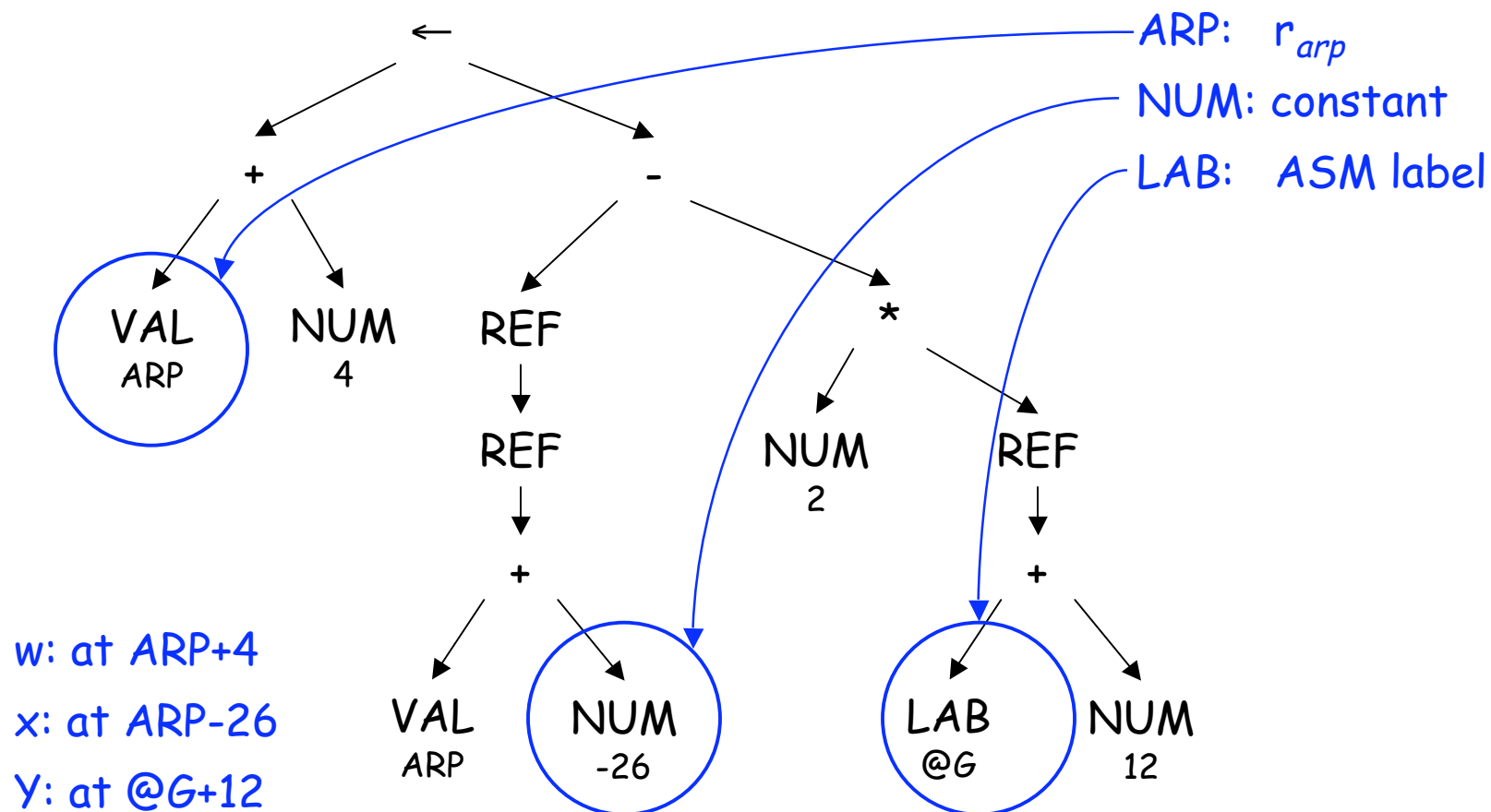


w: at $ARP+4$
x: at $ARP-26$
y: at $@G+12$



The Concept

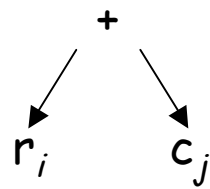
Low-level AST for $w \leftarrow x - 2 * y$



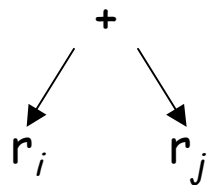
Notation



To describe these trees, we need a concise notation



$+(r_i, c_j)$



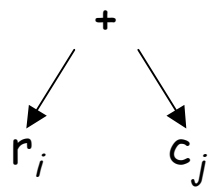
$+(r_i, r_j)$

Linear prefix form

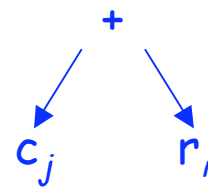
Notation



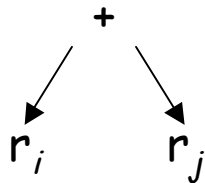
To describe these trees, we need a concise notation



$+(r_i, c_j)$



$+(c_j, r_i)$



$+(r_i, r_j)$

Pattern for commutative variant of $+(r_i, c_j)$

With each tree pattern, we associate a code template and a cost

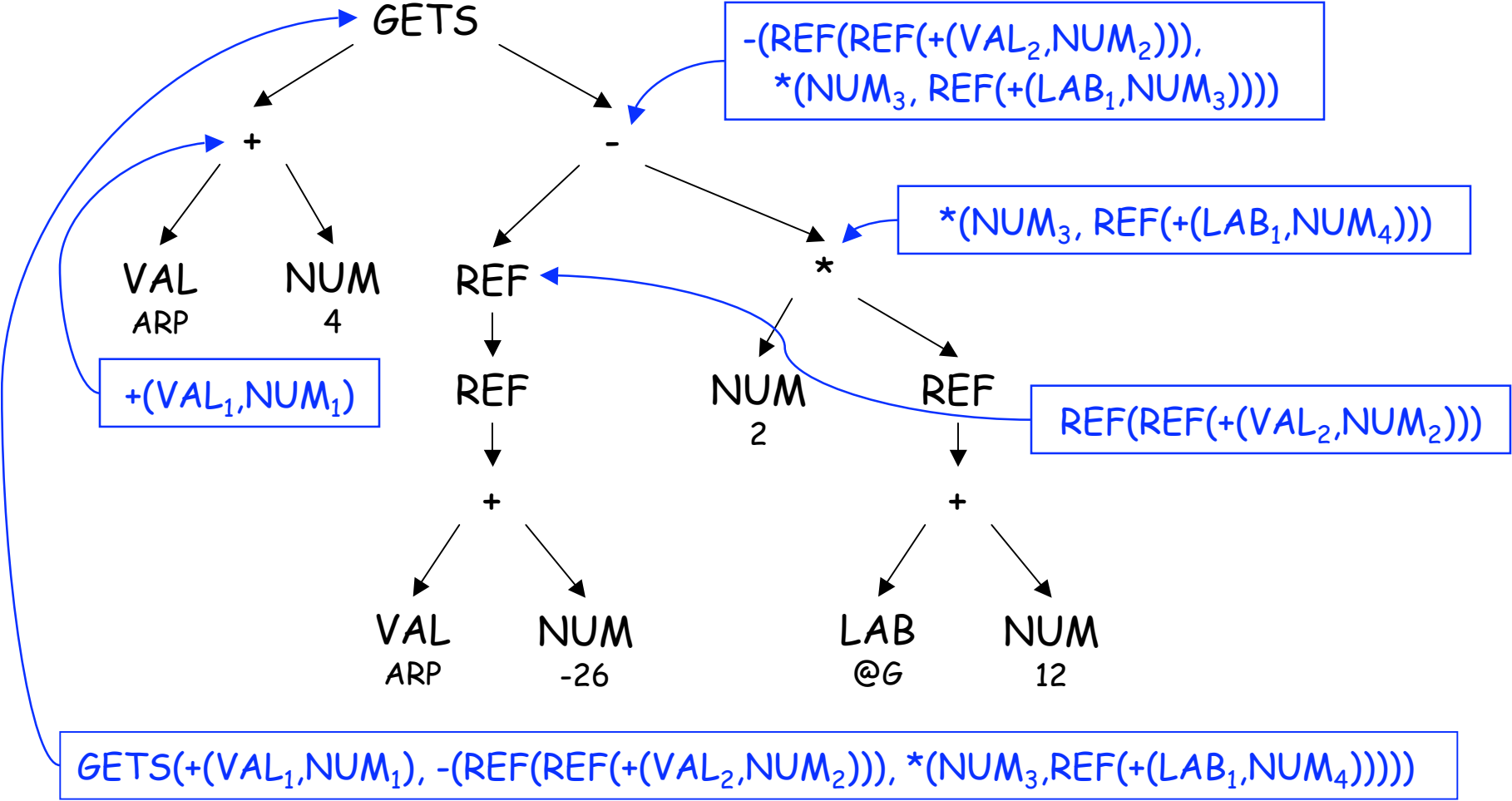
- Template shows how to implement the subtree
- Cost is used to drive process to low-cost code sequences

Subscripts added to create unique names



Notation

The same notation can describe our low-level AST





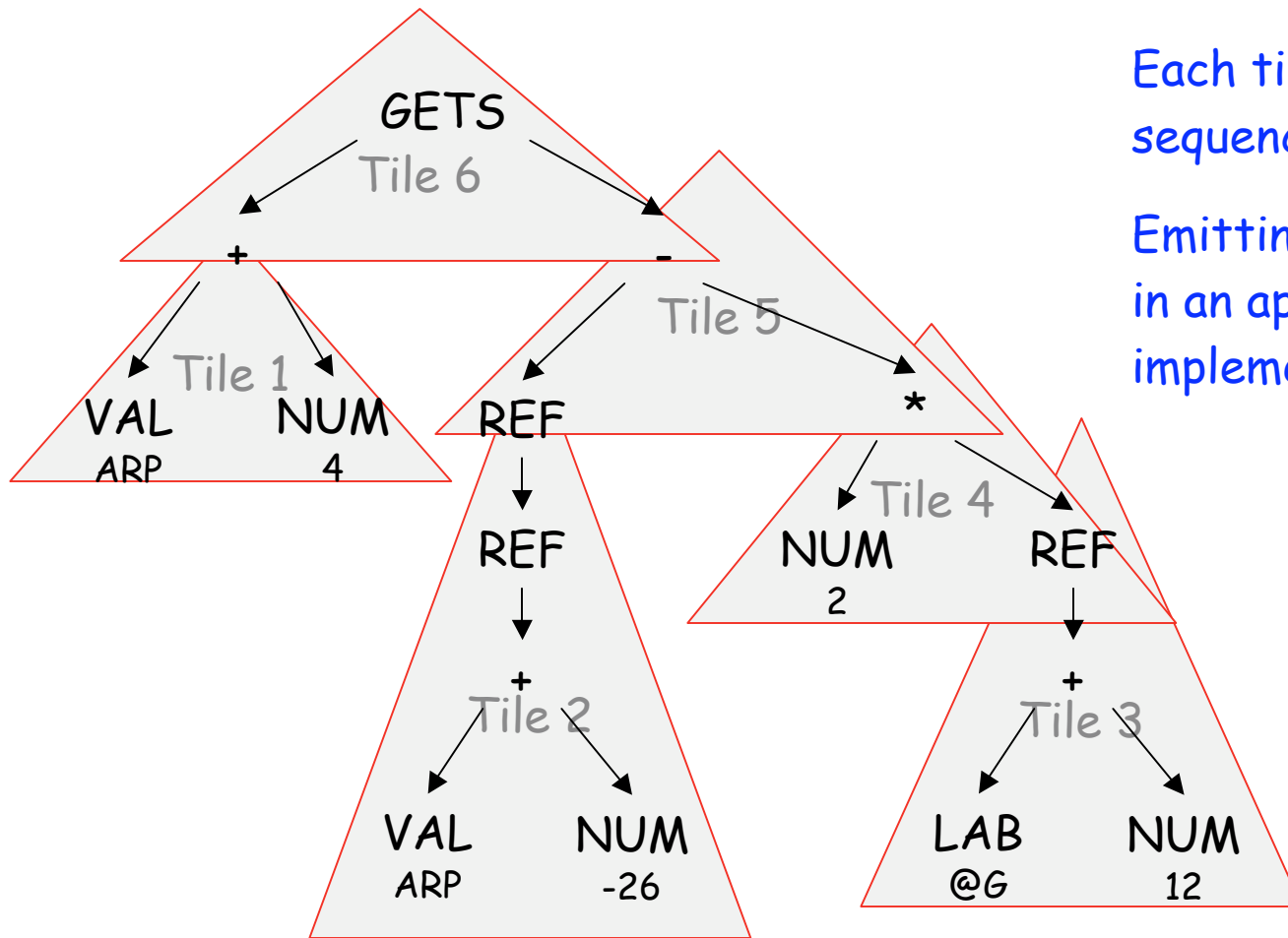
Tree-pattern matching

Goal is to “tile” AST with operation trees

- A tiling is collection of $\langle ast, op \rangle$ pairs
 - ast is a node in the AST
 - op is a pattern tree
 - $\langle ast, op \rangle$ 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
 - $\langle ast, op \rangle \in \text{tiling}$ means that the *root* of ast is also covered by a leaf in another pattern tree in the tiling, unless it is the root
 - Where two pattern trees overlap, they must be compatible (expect the value in the same location)



A Tiling for our Example Tree



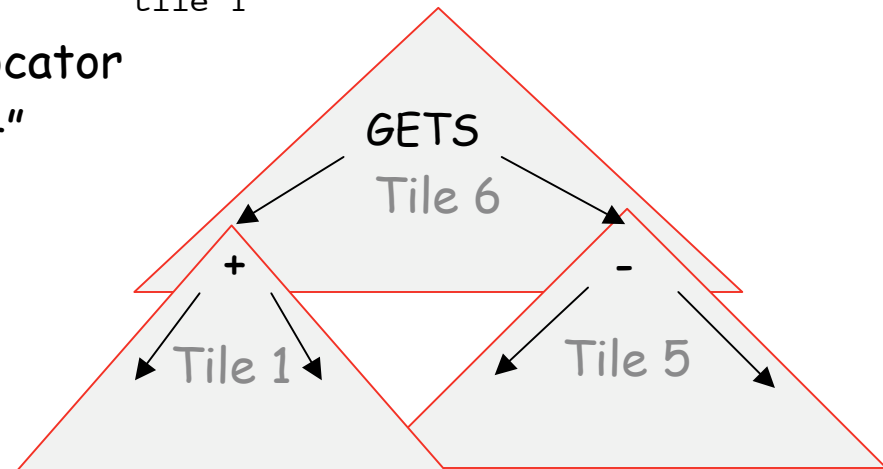
Each tile corresponds to a sequence of operations

Emitting those operations in an appropriate order implements the tree.'



Generating Code from the Tiled Tree

- Postorder treewalk, with node-dependent visitation order
 - Right child of GETS before its left child
 - Might impose "most demanding subtree first" rule ... (Sethi)
- Emit code sequence for tiles, in traversal order
- Tie boundaries together with names
 - Tile 6 uses registers produced by tiles 1 & 5
 - Tile 6 emits "store $r_{\text{tile } 5} \Rightarrow r_{\text{tile } 1}$ "
 - 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
 - Associate a code template with each rule
 - Give the cost of each 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*)



Commutative pair

	Rule	Cost	Template
13	$\text{Reg} \rightarrow + (\text{Reg}_1, \text{Reg}_2)$	1	add $r_1, r_2 \Rightarrow r_{\text{new}}$
14	$\text{Reg} \rightarrow + (\text{Reg}_1, \text{NUM}_2)$	1	addI $r_1, n_2 \Rightarrow r_{\text{new}}$
15	$\text{Reg} \rightarrow + (\text{NUM}_1, \text{Reg}_2)$	1	addI $r_2, n_1 \Rightarrow r_{\text{new}}$
16	$\text{Reg} \rightarrow - (\text{Reg}_1, \text{Reg}_2)$	1	sub $r_1, r_2 \Rightarrow r_{\text{new}}$
17	$\text{Reg} \rightarrow - (\text{Reg}_1, \text{NUM}_2)$	1	subI $r_1, n_2 \Rightarrow r_{\text{new}}$
18	$\text{Reg} \rightarrow - (\text{NUM}_1, \text{Reg}_2)$	1	rsubI $r_2, n_1 \Rightarrow r_{\text{new}}$
19	$\text{Reg} \rightarrow \times (\text{Reg}_1, \text{Reg}_2)$	1	mult $r_1, r_2 \Rightarrow r_{\text{new}}$
20	$\text{Reg} \rightarrow \times (\text{Reg}_1, \text{NUM}_2)$	1	multI $r_1, n_2 \Rightarrow r_{\text{new}}$
21	$\text{Reg} \rightarrow \times (\text{NUM}_1, \text{Reg}_2)$	1	multI $r_2, n_1 \Rightarrow r_{\text{new}}$

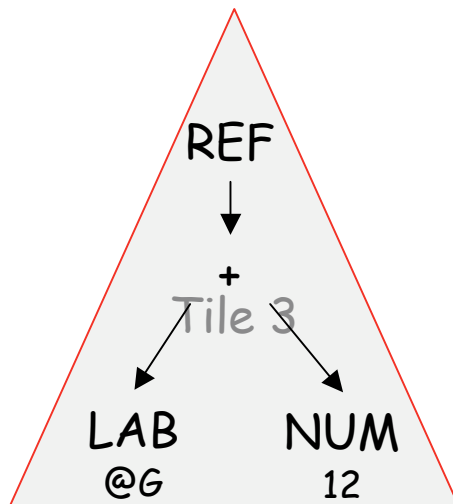
A real set of rules would cover more than signed integers ...



So, What's Hard About This?

Need an algorithm to AST subtrees with the rules

Consider tile 3 in our example



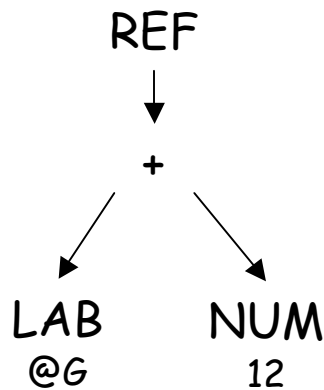


So, What's Hard About This?

Need an algorithm to AST subtrees with the rules

Consider tile 3 in our example

What rules match tile 3?



So, What's Hard About This?

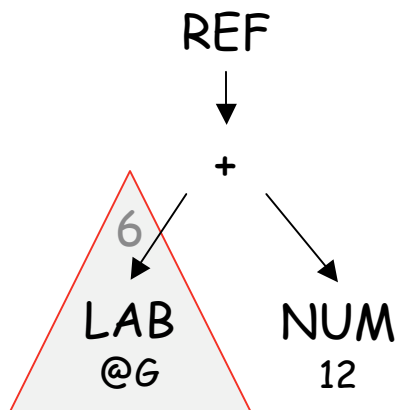


Need an algorithm to match AST subtrees with the rules

Consider tile 3 in our example

What rules match tile 3?

6: $Reg \rightarrow LAB_1$ tiles the lower left node

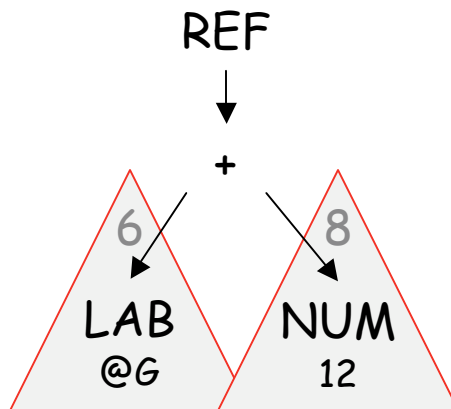




So, What's Hard About This?

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

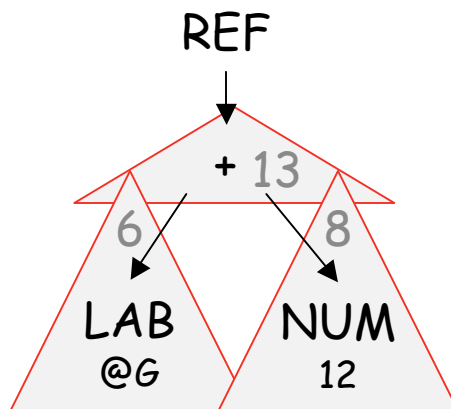
8: $Reg \rightarrow NUM_1$ tiles the bottom right node



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

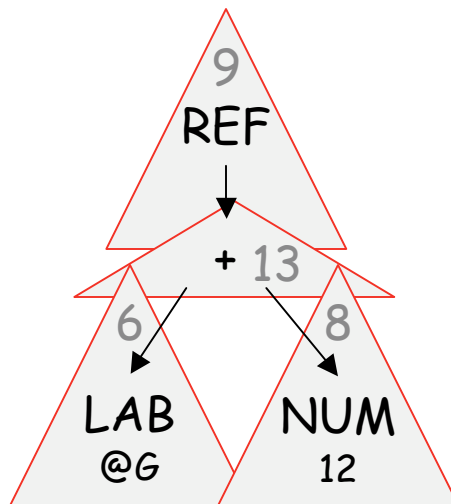
13: $Reg \rightarrow + (Reg_1, Reg_2)$ tiles the + node



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

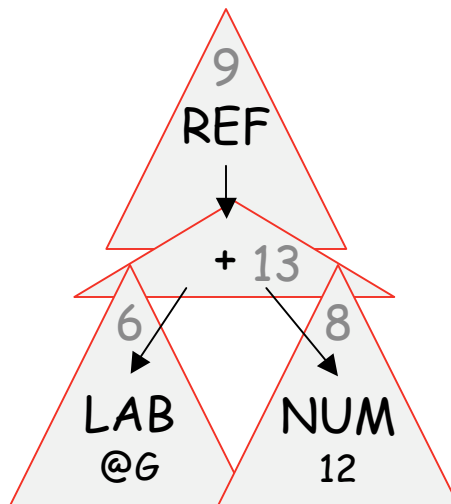
9: $Reg \rightarrow REF(Reg_1)$ tiles the REF



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: $\text{Reg} \rightarrow \text{LAB}_1$ tiles the lower left node

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

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

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

We denote this match as $\langle 6, 8, 13, 9 \rangle$

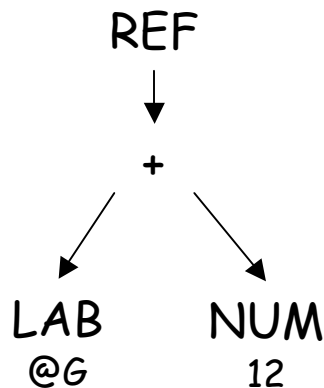
Of course, it implies $\langle 8, 6, 13, 9 \rangle$

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		

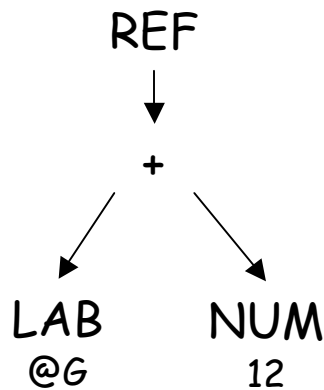
In general, we want the low cost sequence

- We have assumed uniform unit costs
- We favor shorter (lower-cost) sequences
- Accurate costs produce better code

Finding matches



Low Cost Matches



Sequences with Cost of 2	
6: Reg \rightarrow LAB ₁ 11: Reg \rightarrow REF(+ (Reg ₁ , NUM ₂))	loadI @G \Rightarrow r _i loadAI r _i , 12 \Rightarrow r _j
8: Reg \rightarrow NUM ₁ 12: Reg \rightarrow REF(+ (NUM ₁ , Reg ₂))	loadI 12 \Rightarrow r _i loadAI r _i , @G \Rightarrow r _j

These two sequences are equivalent in cost

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

Can encode length restriction on immediate fields into the classification of terminals, at the cost of a little more ambiguity

Tiling the Tree



Still need an algorithm

- Assume each rule implements one operator
 - This is a big assumption because it eliminates rules like 11
$$\text{Reg} \rightarrow \text{REF}(+ (\text{Reg}_1, \text{NUM}_2))$$
 (corresponds to loadAI)
 - A solution: break into two rules
$$\begin{aligned} \text{Reg} &\rightarrow \text{REF}(\text{DUM}) && \text{cost: } 1 \\ \text{DUM} &\rightarrow + (\text{Reg}_1, \text{NUM}_2) && \text{cost: } 0 \end{aligned}$$
- Assume operator takes 0, 1, or 2 operands

Now, ...

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(left 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}

Match binary nodes
against binary rules

Match unary nodes
against unary rules

Handle leaves with
lookup in rule table

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(left 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}

Use same notation to navigate through the rules (pattern trees) and the AST

- *left* and *right* have obvious meanings on AST
- *left* and *right* have similar meanings on the rhs of a rule — interpreted as if on the underlying pattern tree

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(left 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) ← ∅

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(left 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(left 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| × |labels| × |labels|

200 × 1000 × 1000

leads to 200,000,000 entries

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

The Big Picture



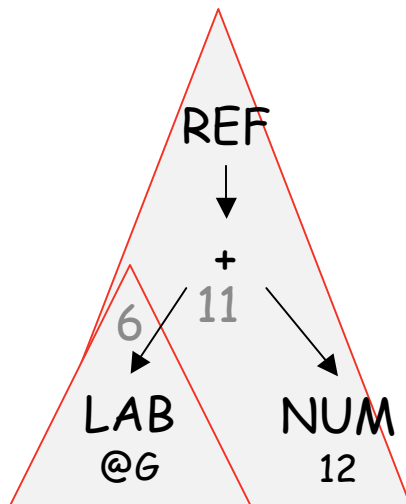
- Tree patterns represent *AST* and *ASM*
- Can use matching algorithms to find low-cost tiling of *AST*
- Can turn a tiling into code using templates for matched rules
- Techniques (& tools) exist to do this efficiently

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



Extra Slides Start Here

Other Sequences for Tile 3



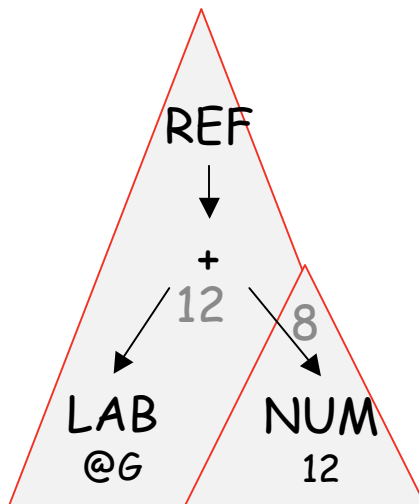
6,11

6: $\text{Reg} \rightarrow \text{LAB}_1$

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

Two operator rule

Other Sequences for Tile 3



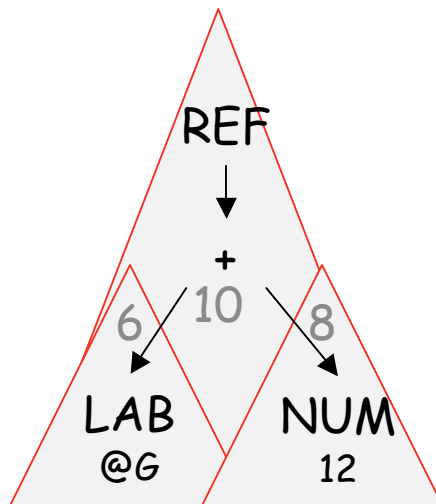
8,12

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

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

Two operator rule

Other Sequences for Tile 3



6,8,10

6: $\text{Reg} \rightarrow \text{LAB}_1$

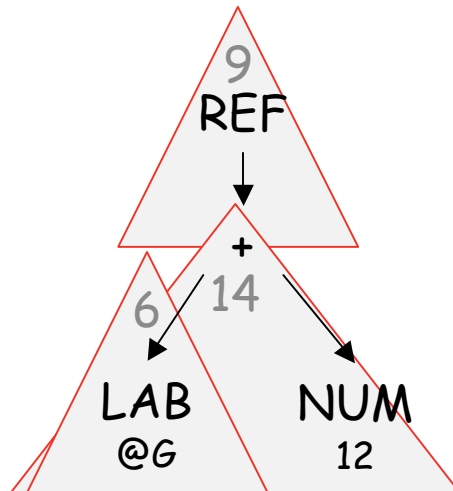
8: $\text{Reg} \rightarrow \text{NUM}_1$

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

Two operator rule

8,6,10 looks the same

Other Sequences for Tile 3



6,14,9

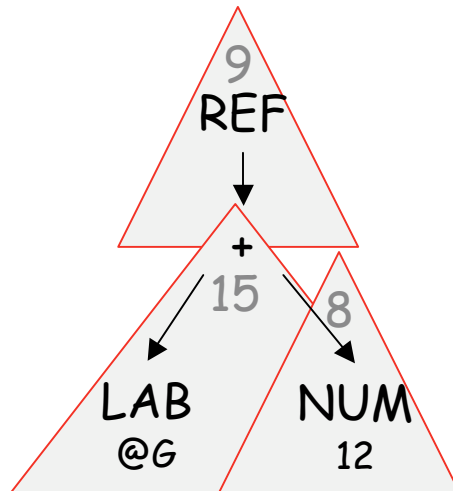
6: $\text{Reg} \rightarrow \text{LAB}_1$

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

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

All single operator rules

Other Sequences for Tile 3



8,15,9

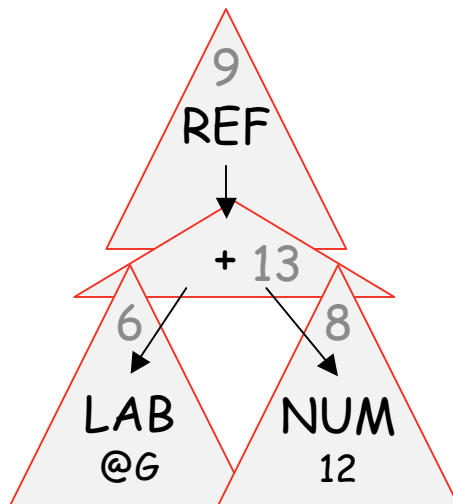
8: $Reg \rightarrow NUM_1$

15: $Reg \rightarrow + (NUM_1, Reg_2)$

9: $Reg \rightarrow REF(Reg_1)$

All single operator rules

Other Sequences for Tile 3



6,8,13,9

6: $\text{Reg} \rightarrow \text{LAB}_1$

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

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

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

All single operator rules

8,6,13,9 looks the same