### CS4212: Compiler Design

## Week 5: LLVMlite Basics of Lexical Analysis

Ilya Sergey

ilya@nus.edu.sg

ilyasergey.net/CS4212/

- Working with data types in LLVM
- LLVMLite Specification
- Overview of HW3
- Lexical Analysis (basics)



#### Announcements

#### • HW3: LLVMlite

- Due: Tuesday, 1 October 2024 at 23:59:59

# - Will be available on Canvas and GitHub on Saturday.



#### Low-Level Virtual Machine (LLVM)

- Open-Source Compiler Infrastructure
  - see llvm.org for full documentation
- Created by Chris Lattner (advised by Vikram Adve) at UIUC
  - LLVM: An infrastructure for Multi-stage Optimization, 2002
  - LLVM: A Compilation Framework for Lifelong Program Analysis and Transformation, 2004
- 2005: Adopted by Apple for XCode 3.1
- Front ends:
  - llvm-gcc (drop-in replacement for gcc)
  - Clang: C, objective C, C++ compiler supported by Apple
  - various languages: Swift, ADA, Scala, Haskell, ...
- Back ends:
  - x86 / Arm / Power / etc.

ported by Apple kell, ...



## LLVM Compiler Infrastructure



#### [Lattner et al.]

- LLVM offers a textual representation of its IR
  - files ending in .ll

factorial64.c

```
#include <stdio.h>
#include <stdint.h>
```

```
int64 t factorial(int64 t n) {
 int64 t acc = 1;
 while (n > 0) {
  acc = acc * n;
  n = n - 1;
 }
 return acc;
```

#### Example LLVM Code

#### factorial-pretty.ll

define @factorial(%n) { %1 = alloca%acc = alloca store %n, %1 store 1, %acc br label %start

start: %3 = 10ad %1%4 = icmp sgt %3, 0br %4, label %then, label %else

then: %6 = load %acc%7 = 10ad %1%8 = mul %6, %7 store %8, %acc %9 = 10ad %1%10 = sub %9, 1 store %10, %1 br label %start

else: %12 = 10ad %accret %12

#### **Real LLVM**

- Decorates values with type information i64 i64\* i1 (boolean)
- Permits numeric identifiers
- Has alignment annotations (padding for some specified number of bytes)
- Keeps track of entry edges for each block: preds = %5, %0

#### factorial.II

; Function Attrs: nounwind ssp define i64 @factorial(i64 %n) #0 { %1 =alloca i64, align 8 %acc = alloca i64, align 8 store i64 %n, i64\* %1, align 8 store i64 1, i64\* %acc, align 8 br label %2 ; <label>:2 ; preds = %5, %0 %3 = load i64\* %1, align 8 %4 = icmp sgt i64 %3, 0br i1 %4, label %5, label %11 ; <label>:5 ; preds = %2 $\%6 = \text{load i64}^*$  %acc, align 8 %7 = load i64\* %1, align 8 %8 = mul nsw i64 %6, %7 store i64 %8, i64\* %acc, align 8 %9 = load i64\* %1, align 8 %10 = sub nsw i64 %9, 1 store i64 %10, i64\* %1, align 8 br label %2 ; preds = %2; <label>:11 %12 = load i64\* %acc, align 8ret i64 %12

### **Example Control-flow Graph**



else: %12 = load %acc ret %12

define @factorial(%n) { %1 = alloca%acc = alloca store %n, %1 store 1, %acc br label %start

start: %3 = 10ad %1%4 = icmp sgt %3, 0br %4, label %then, label %else

then: %6 = load %acc%7 = 10ad %1%8 = mul %6, %7 store %8, %acc %9 = 10ad %1%10 = sub %9, 1 store %10, %1 br label %start

```
else:
 \%12 = \text{load} \%\text{acc}
 ret %12
```

#### LL Basic Blocks and Control-Flow Graphs

- LLVM enforces (some of) the basic block invariants syntactically.
- Representation in OCaml:

type block = {

- - No two blocks have the same label
  - All terminators mention only labels that are defined among the set of basic blocks
  - There is a distinguished, unlabelled, entry block:

insns : (uid \* insn) list; term : (uid \* terminator)

• A control flow graph is represented as a list of labeled basic blocks with these invariants:

type cfg = block \* (lbl \* block) list

## LL Storage Model: Locals

- Several kinds of storage:
  - Local variables (or temporaries): %uid
  - Global declarations (e.g. for string constants): @gid

  - Heap-allocated structures created by external calls (e.g. to malloc)
- Local variables:
  - Defined by the instructions of the form %uid = ...
  - Must satisfy the *single static assignment* (SSA) invariant
  - The value of a %uid remains unchanged throughout its lifetime
  - Analogous to "let %uid = e in …" in OCaml
- Intended to be an abstract version of machine registers.  $\bullet$
- phi functions (<u>https://en.wikipedia.org/wiki/Static\_single-assignment\_form</u>)

– Abstract locations: references to (stack-allocated) storage created by the alloca instruction

• Each %uid appears on the left-hand side of an assignment only once in the entire control flow graph.

Full "SSA" to allow richer use of local variables by taking the control flow into the account



### LL Storage Model: alloca

- The alloca instruction allocates stack space and returns a reference to it.
  - The returned reference is stored in local: %ptr = alloca typ
  - The amount of space allocated is determined by the type
- The contents of the slot are accessed via the load and store instructions:

%acc = alloca i64 store i64 4212, i64\* %acc %x = load i64, i64\* %acc

Gives an abstract version of stack slots

; allocate a storage slot ; store the integer value 4212 ; load the value 4212 into %x



Structured Data

## **Compiling Structured Data**

- Consider C-style structures like those below.
- How do we represent **Point** and **Rect** values?

struct	Point	{ <b>i</b> 1	nt x;	in
struct	Rect	{ s	truct	Po
struct	Rect 1	nk_so	quare	(st
struc	ct Rect	t squ	uare;	
squar	ce.ll =	= squ	lare.	lr
squar	ce.lr.	x +=	len;	
squar	ce.ul.	y +=	len;	
squar	ce.ur.	x +=	len;	
squar	ce.ur.	y +=	len;	
retu	rn squa	are;		

t y; };

int ll, lr, ul, ur };

ruct Point 11, int len) {

= square.ul = square.ur = ll;

## **Representing Structs**

Χ

- Store the data using two contiguous words of memory. lacksquare
- Represent a Point value p as the address of the first word.

struct Rect { struct Point II, Ir, ul, ur };

• Store the data using 8 contiguous words of memory.

- Compiler needs to know the *size* of the struct at compile time to allocate the needed storage space. Compiler needs to know the *shape* of the struct at compile time to index into the structure.

struct Point { int x; int y;};

У

X	lr.y	ul.x	ul.y	ur.x	ur.y
---	------	------	------	------	------



#### Assembly-level Member Access



- Consider: [square.ul.y] = (x86.operand, x86.insns)
- Assume that %rcx holds the base address of square
- Calculate the offset relative to the base pointer of the data: • ul = sizeof(struct Point) + sizeof(struct Point)

$$- y = sizeof(int)$$

• So: [square.ul.y] = (ans, Movq 20(%rcx) ans)

- lr.x lr.y ul.x ul.y ur.x ur.y
- struct Rect { struct Point II, Ir, uI, ur };



#### • How to lay out non-homogeneous structured data?



#### Padding & Alignment

### Copy-in/Copy-out

#### When we do an assignment in C as in:

. . .

struct Rect mk\_square(struct Point II, int elen) { struct Square res; res.Ir = II;

then we copy all of the elements out of the source and put them in the target. Same as doing word-level operations:

struct Rect mk\_square(struct Point II, int elen) { struct Square res; res.lr.x = II.x; res.lr.y = II.x; ...

(which is implemented using a loop in assembly).

• For really large copies, the compiler uses something like memcpy

#### **C** Procedure Calls

- Similarly, when we call a procedure, we copy arguments in, and copy results out.

  - We do the same with scalar values such as integers or doubles.
- Sometimes, this is termed "call-by-value".
  - This is bad terminology.
  - Copy-in/copy-out is more accurate.
- Benefit: locality
- Problem: expensive for large records...
- In C: can opt to pass *pointers* to structs: "call-by-reference" lacksquare

- Caller sets aside extra space in its frame to store results that are bigger than will fit in %rax.

Languages like Java and OCaml always pass non-word-sized objects by reference.

## Representing Data Types







```
void foo() {
 char buf[27];
 buf[0] = 'a';
 buf[25] = 'z';
 buf[26] = 0;
```

- Space is allocated on the stack for buf.
  - need to know size of buf at compile time...
- **buf**[i] is really just (base\_of\_array) + i \* elt\_size

#### Arrays

```
void foo() {
                   char buf[27];
               *(buf) = 'a';
buf[1] = 'b'; *(buf+1) = 'b';
            *(buf+25) = 'z';
                 *(buf+26) = 0;
                  }
```

- Note, without the ability to allocated stack space dynamically (C's alloca function)

#### Multi-Dimensional Arrays

- In C, int M[4][3] yields an array with 4 rows and 3 columns.
- Laid out in *row-major* order:

M[0][0]	M[0][1]	M[0][2]	M[1][0]	M[1][1]	M[1][2]	M[2][0]	
---------	---------	---------	---------	---------	---------	---------	--

• In Fortran, arrays are laid out in *column major order*.

#### M[0][0] M[1][0] M[2][0] M[3][0]

- In ML and Java, there are no multi-dimensional arrays:
   (int array) array is represented as an array of pointers to arrays of ints.
- Why is knowing these memory layout strategies important?

[0] M[0][1] M[1][1] M[2][1] ...

### **Array Bounds Checks**

- ensure that they're in bounds.
  - Compiler generates code to test that the computed offset is legal
- Needs to know the size of the array... where to store it? One answer: Store the size *before* the array contents. \_\_\_\_

- Other possibilities: •
  - Pascal: only permit statically known array sizes (very unwieldy in practice)
  - What about multi-dimensional arrays? \_\_\_\_

• Safe languages (e.g. Java, C#, ML but not C, C++) check array indices to

A[2] A[3] A[4] A[5] A[6]

#### Array Bounds Checks (Implementation)

To read a value from the array **arr**[i]:

	movq -8(%rax) %rdx	// loa
	cmpq %rdx %rcx	// со
	jIok	// ju
	callqerr_oob	// tes
ok:		
	movq (%rax, %rcx, 8) dest	// dc

- Clearly more expensive: adds move, comparison & jump
  - More memory traffic
- These overheads are particularly bad in an inner loop
- Compiler optimisations can help remove the overhead
  - e.g. In a for loop, if bound on index is known, only do the test once

Example: Assume %rax holds the base pointer (arr) and %rcx holds the array index i.

ad size into rdx mpare index to bound mp if  $0 \le i \le size$ st failed, call the error handler

o the load from the array access

– Hardware can improve performance: executing instructions in parallel, branch prediction

## **C-style Strings**

- A string constant "foo" is represented as global data: \_string42: 102 111 111 0
- C uses null-terminated strings
- allows all copies of the same string to be shared.
- Rookie mistake (in C): write to a string constant.

- Instead, must allocate space on the heap:

char \*p = (char \*)malloc(4 \* sizeof(char)); strncpy(p, "foo", 4); /\* include the null byte \*/ p[0] = 'b';

• Strings are usually placed in the *text* segment so they are *read only*.

char \*p = "foo"; p[0] = 'b';

Attempting to modify the string literal is *undefined behaviour*.





#### **C-style Enumerations / ML-style datatypes**

- In C:
- In OCaml: type day = Sun | Mon | Tue | Wed | Thu | Fri | Sat ullet
- C lets programmers choose the tags
- $\bullet$
- Representation: a foo value is a pointer to a pair: (tag, data)
- Example: tag(Bar) = 0, tag(Baz) = 1 $\llbracket \text{let f} = \text{Bar}(3) \rrbracket =$

[[let g = Baz(4, f)]] =

enum Day {sun, mon, tue, wed, thu, fri, sat} today;

Associate an integer *tag* with each case: sun = 0, mon = 1, ...

OCaml datatypes can also carry data: type foo = Bar of int | Baz of int \* foo



### Switch Compilation

• Consider the C statement:

swi	.tch (	e) {		
	case	sun:	s1;	bre
	case	mon:	s2;	bre
	•••			
	case	sat:	s3;	bre
}				

- How to compile this?
  - What happens if some of the break statements are omitted? (Control falls through to the next branch.)

ak;

ak;

ak;

## Cascading ifs and Jumps

 $[switch(e) \{case tag1: s1; case tag2 s2; ...\}] =$ 

- Each \$tag1...\$tagN is just a constant int tag value.
- Note: [break;]
  (within the switch branches) is:

br %merge

```
%tag = [[e]];
br label %l1
```

```
I1: %cmp1 = icmp eq %tag, $tag1
    br %cmp1 label %b1, label %l2
b1: [s1]
```

```
br label %l2
```

```
l2: %cmp2 = icmp eq %tag, $tag2
    br %cmp2 label %b2, label %l3
b2: [s2]
```

```
br label %l3
```

```
IN: %cmpN = icmp eq %tag, $tagN
    br %cmpN label %bN, label %merge
bN: [sN]
    br label %merge
```

merge:



### **Alternatives for Switch Compilation**

- Nested if-then-else works OK in practice if # of branches is small - (e.g. < 16 or so).
- For more branches, use better data structures to organise the jumps:
  - Create a table of pairs (v1, branch\_label) and loop through
  - Or, do binary search rather than linear search
  - Or, use a hash table rather than binary search
- One common case: the tags are dense in some range [min...max]
  - Let N = max min

  - Create a branch table Branches[N] where Branches[i] = branch\_label for tag i. Compute tag = [e] and then do an *indirect jump*: J Branches[tag]
- Common to use heuristics to combine these techniques.  $\bullet$

## **ML-style Pattern Matching**

- ML-style match statements are like C's switch statements except:  $\bullet$ 
  - Patterns can bind variables
  - Patterns can nest —

- Compilation strategy:  $\bullet$ 
  - "Flatten" nested patterns into matches against one constructor at a time.
  - Compile the match against the tags of the datatype \_\_\_\_ as for C-style switches.
  - \_\_\_\_
- - Many of these transformations can be done at the AST level —



Code for each branch additionally must copy data from [e] to the variables bound in the patterns.

• There are many opportunities for optimisations, many papers about "pattern-match compilation"



Good place for a break



#### Datatypes in LLVM IR

#### Structured Data in LLVM

• LLVM's IR is uses types to describe the structure of data.

- <#elts> is an integer constant >= 0
- Structure types can be named at the top level:

$$T1 = type$$

• Such structure types can be recursive

N-bit integers arrays function types structures pointers named (identified) type

ction Types *return, argument types* 

e { $t_1$ ,  $t_2$ , ...,  $t_n$ }

## **Example LL Types**

- A static array of 4212 integers:
- A two-dimensional array of integers: [ 3 x [ 4 x i64 ] ]
- Structure for representing dynamically-allocated arrays with their length: { i64 , [0 x i64] }
  - There is no array-bounds check; the static type information is only used for calculating pointer offsets.
- C-style linked lists (declared at the top level): %Node = type { i64, %Node\*}
- Structs from the C program shown earlier: %Rect = type { %Point, %Point, %Point, %Point } %Point = type { i64, i64 }

[ 4212 x i64 ]
#### getelementptr

- LLVM provides the getelementptr instruction to compute pointer values
  - Given a pointer and a "path" through the structured data pointed to by that pointer, getelementptr computes an address

  - It is a "type indexed" operation, since the size computations depend on the type

insn ::= ... getelementptr t\* %val, t1 idx1, t2 idx2 ,...

Example: access the x component of the first point of a rectangle:

%tmp1 = getelementptr %Rect\* %square, i32 0, i32 0 %tmp2 = getelementptr %Point\* %tmp1, i32 0, i32 0

• The first is i32 0 a "step through" the pointer to, e.g., %square, with offset 0.

See "Why is the extra 0 index required?": https://llvm.org/docs/GetElementPtr.html#why-is-the-extra-0-index-required

– This is the abstract analog of the X86 LEA (load effective address). It **does not** access memory.

### **GEP Example\***



Final answer: ADDR + size ty(%ST) + size\_ty(%RT) + size\_ty(i32) + size ty(i32) + 5\*20\*size ty(i32) + 13\*size ty(i32)

\*adapted from the LLVM documentation: see <u>http://llvm.org/docs/LangRef.html#getelementptr-instruction</u>

#### getelementptr

- GEP *never* dereferences the address it's calculating:
  - GEP only produces pointers by doing arithmetic
  - It doesn't actually traverse the links of a data structure
- To index into a deeply nested structure, one has to "follow the pointer" by loading from the computed pointer

### **Compiling Data Structures via LLVM**

- Translate high level language types into an LLVM representation type. 1. For some languages (e.g. C) this process is straightforward — The translation simply uses platform-specific alignment and padding  $\bullet$
- - For other languages, (e.g. OO languages) there might be a fairly complex elaboration.
    - e.g. for OCaml, arrays types might be translated to pointers to length-indexed structs. lacksquare $[int array] = \{i32, [0 \times i32]\}*$
- Translate accesses of the data into getelementptr operations: 2.
  - e.g. for OCaml array size access: \_\_\_\_ [length a] =
    - $\$1 = getelementptr \{i32, [0 x i32]\} * \$a, i32 0, i32 0$

- What if the LLVM IR's type system isn't expressive enough?

  - e.g. if the source language has polymorphic/generic types
- LLVM IR provides a bitcast instruction
  - (segmentation faults, or silent memory corruption)

%rect2 = type { i64, i64 } %rect3 = type { i64, i64, i define @foo() { %1 = alloca %rect3 ; %2 = bitcast %rect3\* %1 t %3 = getelementptr %rect2

### Type Casting

– e.g. if the source language has subtyping, perhaps due to inheritance

– This is a form of (potentially) unsafe cast. Misuse can cause serious bugs

.64	: }	; ;	tw th	o-fi ree-	leld fie	re ld	ecoro	d ord	
al	loca	ate a	a th	ree-	-fie	ld	reco	ord	
0	%rec	:t2*		; 58	afe	ca	st		
*	%2 <b>,</b>	i32	0,	i32	1	; 3	allo	wed	

https://ilyasergey.net/CS4212/hw03-llvmlite-spec.html

## LLVMlite Specification

#### LLVMlite features

- A C-like "weak type system" to statically rule out some malformed programs.
- A variety of different kinds of integer values, pointers, function pointers, and structured data including strings, arrays, and structs.
- Top-level mutually-recursive function definitions and function calls as primitives.
- An infinite number of "locals" (also known as "pseudo-registers", "SSA variables", or  $\bullet$ "temporaries") to hold intermediate results of computations.
- An abstract memory model that doesn't constrain the layout of data in memory.  $\bullet$
- Dynamically allocated memory associated with a function invocation (in C, the stack).  $\bullet$  $\bullet$
- Static and dynamically (heap) allocated structured data.
- A control-flow graph representation of function bodies.



```
define i64 @fac(i64 %n) {
 %1 = icmp sle i64 %n, 0
  br i1 %1, label %ret, label %rec
ret:
  ret i64 1
rec:
  %2 = sub i64 %n, 1
 %3 = call i64 @fac(i64 %2)
 %4 = mul i64 %n, %3
  ret i64 %4
}
define i64 @main() {
 %1 = call i64 @fac(i64 6)
  ret i64 %1
J
```

## Example

;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;;	<pre>(1) (2) (3) (4) (5) (6)</pre>	function definition, argument prefixed with % signed comparison, result assigned to %1 "terminator", marks the end of the block label, indicates the beginning of the new block return the result (1) another block subtract 1 from %n_name result %2
;	(7)	call function @fac, assign the result for %3
;	(8)	return result
;	(9)	call @fac with the argument 6

# LLVMlite types

Concrete Syntax	Kind	Description
void	void	Indicates the instruction does not return a usable value.
i1, i64	simple	1-bit (boolean) and 64-bit integer values.
T*	simple	Pointer that can be dereferenced if its target is compatible with T
i8*	simple	Pointer to the first character in a null-terminated array of bytes. Note: 18* is a valid type, but just 18 is not. LLVMlite programs do not operate over byte-sized integer values.
F*	simple	Function pointer
S(S1,, SN)	function	A function from S1,, SN to S
void(S1,, SN)	function	A function from S1,, SN to void
{ T1,, TN }	aggregate	Tuple of values of types T1,, TN
[N×T]	aggregate	Exactly N values of type T
%NAME	*	Abbreviation defined by a top-level named type definition

- Simple types appear on stack and as arguments to functions •
- $\bullet$
- One can define abbreviations for types: %IDENT = type T

Aggregate types that may only appear in global and heap-allocated data

#### **Global Definitions**

#### @IDENT = global T G

Concrete Syntax	Туре	
null	T*	•
[0-9]+	i64	(
@IDENT	T*	(
c"[A−z]*\00"	[Nxi8]	-
[ T G1,, T GN ]	[N×T]	
{ T1 G1,, TN GN }	{T1,,TN}	
bitcast (T1* G1 to T2*)	T2*	

@foo = global i64 42
@bar = global i64\* @foo
@baz = global i64\*\* @bar

#### Description

The null pointer constant.

64-bit integer literal.

Global identifier. The type is always a pointer of the type associated with the global definition.

String literal. The size of the array N should be the length of the string in bytes, including the null terminator  $\sqrt{00}$ .

Array literal.

Struct literal.

Bitcast.

### Operands of functions

Concrete Syntax	Туре	Description
null	T*	The null pointer constant
[0-9]+	i64	64-bit integer literal
@IDENT	<b>T</b> *	Global identifier. The type can always be determined from the global definitions and is always a pointer
%IDENT	S	Local identifier: can only name values of simple type. The type determined by an local definition of <b>%IDENT</b> in scope

### Types of instructions



- Let's discuss the meaning of these types
- (see the specification)

	Operand $\rightarrow$ Result Types
	i64 × i64 → i64
	$- \rightarrow S*$
	S* → S
	S x S* → void
	$S \times S \rightarrow i1$
	$S1(S2,, SN) * \times S2 \times \times SN \rightarrow S1$
	$void(S2,, SN) * \times S2 \times \times SN \rightarrow void$
PN	T1* x i64 x x i64 -> GEPTY(T1, OP1,, OPN)*
	T1* → T2*

• The getelementptr instruction has some additional well-formedness requirements T

# **GEP** Type

GEPTY : T -> operand list -> T GEPTY T operand::path' = GEPTY' T path' GEPTY': T -> operand list -> T GEPTY'T GEPTY' [ \_ x T ] operand::path' = GEPTY' T path'

- GEPTY is a partial function.  $\bullet$
- When GEPTY is not defined, the corresponding instruction is malformed.
- This happens when, for example:
  - The list of index operands provided is empty
  - An operand used to index a struct is not a constant
  - The type is not an aggregate and the list of indices is not empty

```
|| = T
GEPTY' { T1, ..., TN } (Const m)::path' = GEPTY' Tm path' when m <= N
```

#### Notes on GEP

```
%struct = type { i64, [5 x i64], i64}
@gbl = global %struct {i64 1,
   [5 x i64] [i64 2, i64 3, i64 4, i64 5, i64 6], i64 7}
define void @foo() {
 %1 = getelementptr %struct* @gbl, i32 0, i32 0
  ...
```

- LLVMlite ignores the i32 annotation and treats these as i64 values lacksquare

  - we assume the arguments of getelementptralways fall in the range [0, Int32.max int].

• Real LLVM requires that constants appearing in getelementptr be declared with type i32:

we keep the i32 annotation in the syntax to retain compatibility with the clang compiler

#### Blocks, CFGs, and Function Definitions

Concrete Syntax	
ret void	
ret S OP	
br label %LAB	
br i1 OP, label %LAB1,	labe

- The body of a function is represented by a control flow graph (CFG).
- The full syntax of a function definition:

• A block is just a sequence of instructions followed by a terminator

	Operand $\rightarrow$ Result Types
	$- \rightarrow -$
	S → -
	$- \rightarrow -$
l %LAB2	i1 → -

A CFG consists of a distinguished entry block and a sequence blocks of prefixed with a label.

define [S|void] @IDENT(S1 OP, ..., SN OP) { BLOCK (LAB: BLOCK)...}



#### LLVMlite Semantics

- Like for X86lite, we define the semantics of LLVMlite by describing the execution of an *abstract machine*.
- LLVMlite machine explicitly differentiates between stack, heap, code, and global memory (X86 was treating all of those uniformly).
  A definitional (reference) interpreter for LLVMlite is provided in HW3:
- A definitional (reference) interpr check **llinterp.ml**
- If you have a question about a detail of the semantics, you can simply run a program through the interpreter!

### Memory Model

• The memory state of the LLVMlite machine is represented by a mapping between **block** identifiers and memory values.

> / \ L node / | | \ L L L\* L

- To identify the leaf marked \* we provide the indices 0, 1, 2 along with the identifier **bid1**.

We will refer to a top-level memory value that is not a subtree of another as a **memory block**.

Even simple values, such as a single global i64 will be represented using a node with one leaf. • This means that we're selecting the 2nd child of the 1st child of the 0th child of the root node.

### Memory Model

/ node / | | \ L L L\* L

- The simple values include:  $\bullet$ 
  - 1-bit (boolean) and 64-bit 2's complement signed integers

  - Pointers to a subtree of a particular memory block containing a block identifier and path • A special **undef** value that represents an unusable value

#### Interpreter

- interp\_call takes
- the global identifier of a function in an LLVMlite program,
- a list of (simple) values to serve as arguments, and
- an initial memory state; and
- returns the memory state after the function call has completed and the return value.
- interp\_cfg does most of the work. It takes
  - a control-flow graph,
  - an initial locals map, and
  - a memory state; and
  - evaluates the cfg, returning the new memory state and the return value of the function body.

#### **Some Instructions**

	. <u> </u>
%L = alloca S	Allo poi the frar
%L = load S* OP	OP refe Up poi me its op the sim the
store S OP1, S* OP2	Upo to t the typ cras

(see implementation)

ocate a slot in the current stack frame and return a inter to it. This involves adding a subtree of undef to root node of the memory block representing the me at the next available index.

must be a pointer or **undef**. Find the value erenced by the pointer in the current memory state. date locals( %L ) with the result. If OP is not a valid nter, either because it evaluates to **undef**, no mory value is associated with its block identifier or path does not identify a valid subtree, then the eration raises an error and the machine crashes. If pointer is valid, but the value in memory is not a ple value of type S, the operation raises an error and machine crashes.

date the memory state by setting the target of OP2 the value of OP1. If OP2 is not a valid pointer, or if target of OP2 is not a simple value in memory of be S, the operation raises an error and the machine shes.

# **Some Instructions (c'd)** (see implementation)

%L = call S1 OP1(S2 OP2,, SN (	OPN)	Eva invo cur fun typ ope Upo con
call void OP1(S2 OP2,, SN OPN)	)	The wit

aluate all of the operands and use them to recursively oke the interpreter through **interp\_call** with the rrent memory state. If OP1 does not evaluate to a action pointer that identifies a function with return be S1 and argument types S2, ..., SN, then the eration raises an error and the machine crashes. date the local (%L) to the result of **interp\_call** and ntinue with the return memory state.

e same as a non-void call, but no locals are updated th the returned value.

#### Some Instructions (c'd) (see implementation)

%L =	gete	eleme	ent	ptr	T1*	0P	91,			Cre
i64 (	)P2,		,	i64	<b>OPN</b>					OP
										the
										lf tł
										me
										loca
										loca
										sec

ate a new pointer by *adding* the first index operand to the last index of the pointer value of OP1 and concatenating the remaining indices onto the path. e target of the resulting pointer is not a valid nory value *compatible* with the type %L, then update Is( %L ) with the **undef** value. Otherwise, update Is( %L ) with the new pointer. See the following tion for a more detailed explanation.

```
%t1 = type { A, B, C }
%t2 = type [ 2 x %t1 ]
@pn1 = global %t2 [ {a0, b0, c0}, {a1, b1, c1} ]
; Memory:
  { ... bid0 -> root ... }
                 n1
n3
             n2
a0 b0 c0 a1 b1 c1
%pn2 = getelementptr %t2* pn1, i32 0, i32 0
                                               ; %t1* -> n2
                                               ; B* -> b1
%pb1 = getelementptr %t1* pn2, i32 1, i32 1
```

- Check out **effective** tag in the interpreter code.
- Some examples in llprograms: gep3.ll, gep5.ll, gep6.ll

# **GEP Indexing**

- Start with the pointer pn1 = (bid0, 0) pointing to n1.  $\bullet$
- The first GEP instruction above will compute the pointer (bid0, 0, 0), by first adding 0 to the last index of **pn1** and then concatenating the rest of the indices to the end of the path.
- The next GEP instruction will compute the pointer  $\bullet$ (**bid0**, **0**, **1**, **1**), which points to **b1**.
- Why?  $\bullet$
- Indexing into a sibling (rather than a child) of a node  $\bullet$ using GEP with a non-zero first index is only legal if sibling nodes are allocated as *part of an array*.
- In our example, n1 was allocated as [ 2 x %t1 ], so this  $\bullet$ is the case.

### Compiling LLVMlite to X86

### **Compiling LLVMlite Types to X86**

- [[i1], [[i64]], [[t\*]] = quad word (8 bytes, 8-byte aligned)
- raw i8 values are not allowed (they must be manipulated via i8\*)
- array and struct types are laid out sequentially in memory
- getelementptr computations must be relative to the LLVMlite size definitions
  i.e. [[i1]] = quad (quite wasteful!)

# **Compiling LLVM locals**

- Option 1: lacksquare
  - Map each %uid to a x86 register
  - Efficient!
  - Difficult to do effectively: many %uid values, only 16 registers
  - We will see how to do this later in the semester
- Option 2:
  - Map each %uid to a stack-allocated space
  - Less efficient!
  - Simple to implement
- For HW3 we will follow Option 2

#### • How do we manage storage for each %uid defined by an LLVM instruction?

#### **Other LLVMlite Features**

- Globals lacksquare
  - must use %rip relative addressing
- Calls
  - Follow x64 AMD ABI calling conventions
  - Should interoperate with C programs
- getelementptr
  - trickiest part

- See HW3 description and README.md
- Main definitions: ll.ml lacksquare
- Compiler in the pipeline: driver.ml and process\_ll\_file.
- Using main.native
- Compiling with clang

#### **Tour of HW3**

#### Lexical analysis, tokens, regular expressions, automata



#### **Compilation in a Nutshell**

Source Code (Character stream)







### Today: Lexing

Source Code (Character stream) if (b == 0) { a = 1; }







#### First Step: Lexical Analysis

• Change the character stream "if (b == 0) a = 0;" into tokens:

EQ; Int(0); SEMI; RBRACE

- Token: data type that represents indivisible "chunks" of text:
  - Identifiers: a y11 elsex \_100
  - Keywords: if else while
  - Integers: 2 200 –500 5L
  - Floating point: 2.0 .02 1e5
  - Symbols: + \* ` { } ( ) ++ << >> >>>
  - Strings: "x" "He said, "Are you?""
  - Comments: (\* CS4212: Project 1 ... \*) /\* foo \*/
- Often delimited by *whitespace* (' ', \t, etc.) – In some languages (e.g. Python or Haskell) whitespace is significant

0 a  $\blacksquare$ •

IF; LPAREN; Ident("b"); EQEQ; Int(0); RPAREN; LBRACE; Ident("a");



See handlex.ml

https://github.com/cs4212/week-05-lexing

#### Demo: Handlex

How hard can it be?

- How hard can it be?
  - Tedious and painful!
- Problems:
  - Precisely define tokens
  - Matching tokens simultaneously
  - Reading too much input (need look ahead)
  - Error handling
  - Hard to compose/interleave tokeniser code
  - Hard to maintain


# A Principled Solution to Lexing

## **Regular Expressions**

- Regular expressions precisely describe sets of strings.
- A regular expression R has one of the following forms:

- 8	Epsilon stands for the
- 'a'	An ordinary characte
- R <sub>1</sub>   R <sub>2</sub>	Alternatives, stands for
$- R_1 R_2$	Concatenation, stanc
– R*	Kleene star, stands fo

Useful extensions: lacksquare

	"foo"	Strings, equivalent to
_	R+	One or more repetition
_	R?	Zero or one occurren
_	['a'-'z']	One of a or b or c or
_	[^'0'-'9']	Any character except
_	R as x	Name the string mate

e empty string er stands for itself or choice of  $R_1$  or  $R_2$ ds for  $R_1$  followed by  $R_2$ or zero or more repetitions of R

'f''o''o' ons of R, equivalent to RR\* nces of R, equivalent to  $(\varepsilon | R)$  $\dots z$ , equivalent to (a|b|...|z)0 through 9 ched by  $\mathbf{R}$  as  $\mathbf{x}$ 

## Example Regular Expressions

- Recognise the keyword "if": "if"
- Recognise a digit: ['0'-'9']
- Recognise an integer literal: '-'?['0'-'9']+
- Recognise an identifier: (['a'-'z'] | ['A'-'Z'])(['0'-'9'] |' '| ['a'-'z'] | ['A'-'Z'])\*
- In practice, it's useful to be able to *name* regular expressions:
- let lowercase = ['a'-'z']
- let uppercase = ['A' 'Z']
- let character = uppercase | lowercase

## How to Match?

- Consider the input string: ifx = 0– Could lex as: or as: if X 0 =
- Regular expressions alone are *ambiguous*, need a rule to choose between the options above  $\bullet$
- Most languages choose "longest match"  $\bullet$ 
  - So the 2<sup>nd</sup> option above will be picked
  - Note that only the first option is "correct" for parsing purposes
- Conflicts: arise due to two tokens whose regular expressions have a shared prefix lacksquareTies broken by giving some matches higher priority —
- - Example: keywords have priority over identifiers
  - Usually specified by order the rules appear in the lex input file



## Lexer Generators

- Reads a list of regular expressions:  $R_1, \dots, R_n$ , one per token.
- Each token has an attached "action"  $A_i$ (just a piece of code to run when the regular expression is matched)



- Generates scanning code that: •
  - Decides whether the input is of the form  $(R_1 | ... | R_n) *$
  - Whenever the scanner matches a (longest) token, it runs the associated action 2.

```
{ Int (Int32.of string (lexeme lexbuf)) }
{ PLUS }
{ Ident (lexeme lexbuf) }
{ token lexbuf }
                    actions
```



#### lexlex.mll