

# Compilers & Domain-Specific Languages

Roland Leißa



# **Introduction**

---

## Myth: Computers are fast enough



Monsters University, 2013, Pixar

## Myth: Computers are fast enough

- Rendering a single frame: >24 hours



Monsters University, 2013, Pixar

## Myth: Computers are fast enough



Monsters University, 2013, Pixar

- Rendering a single frame: >24 hours
- Rendering whole movie on a supercomputer:
  - 2000 machines
  - 24,000 cores
  - 2 years to render

## Myth: Computers are fast enough



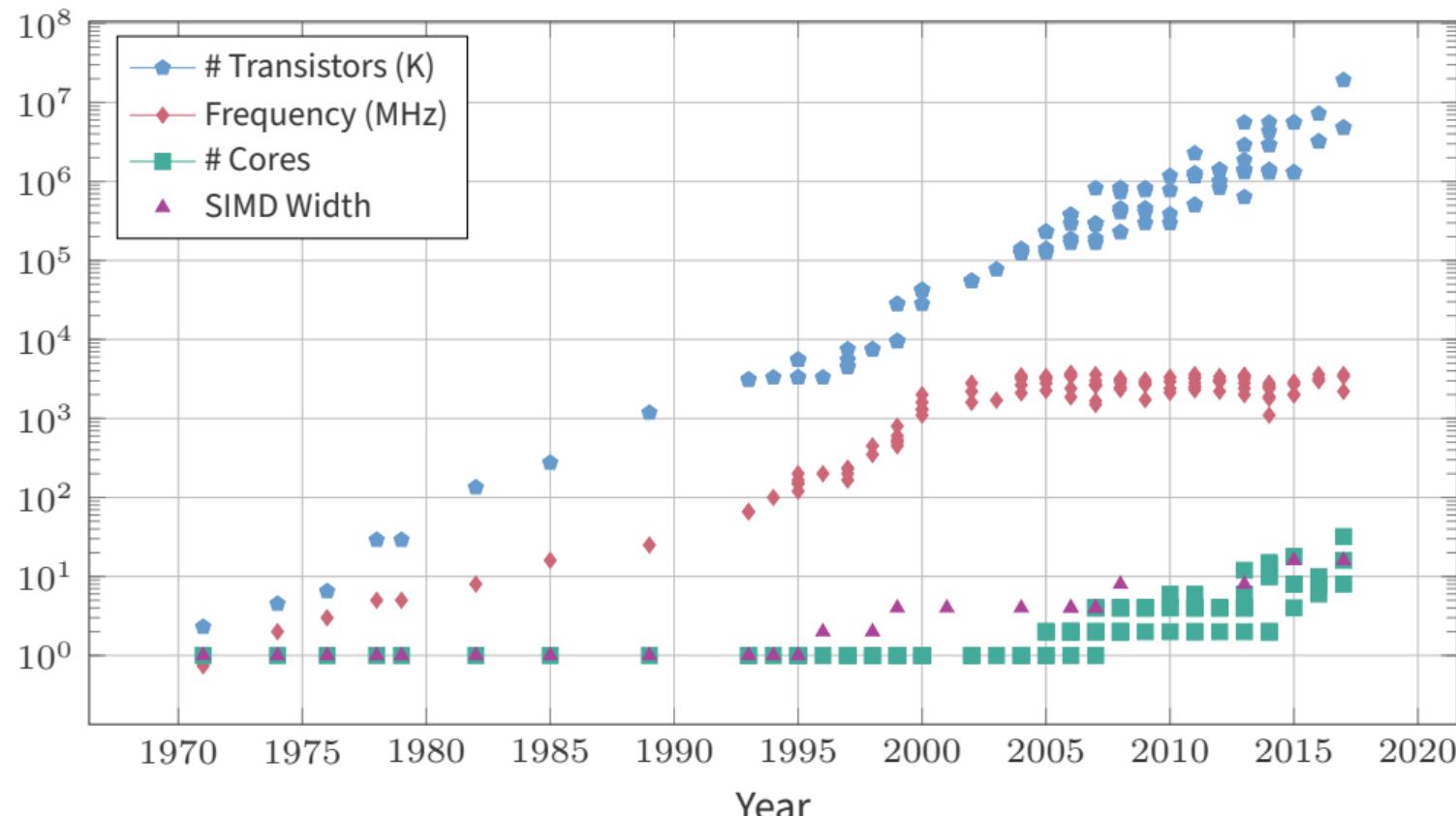
Monsters University, 2013, Pixar

- Rendering a single frame: >24 hours
- Rendering whole movie on a supercomputer:
  - 2000 machines
  - 24,000 cores
  - 2 years to render
- Doing this in real time?

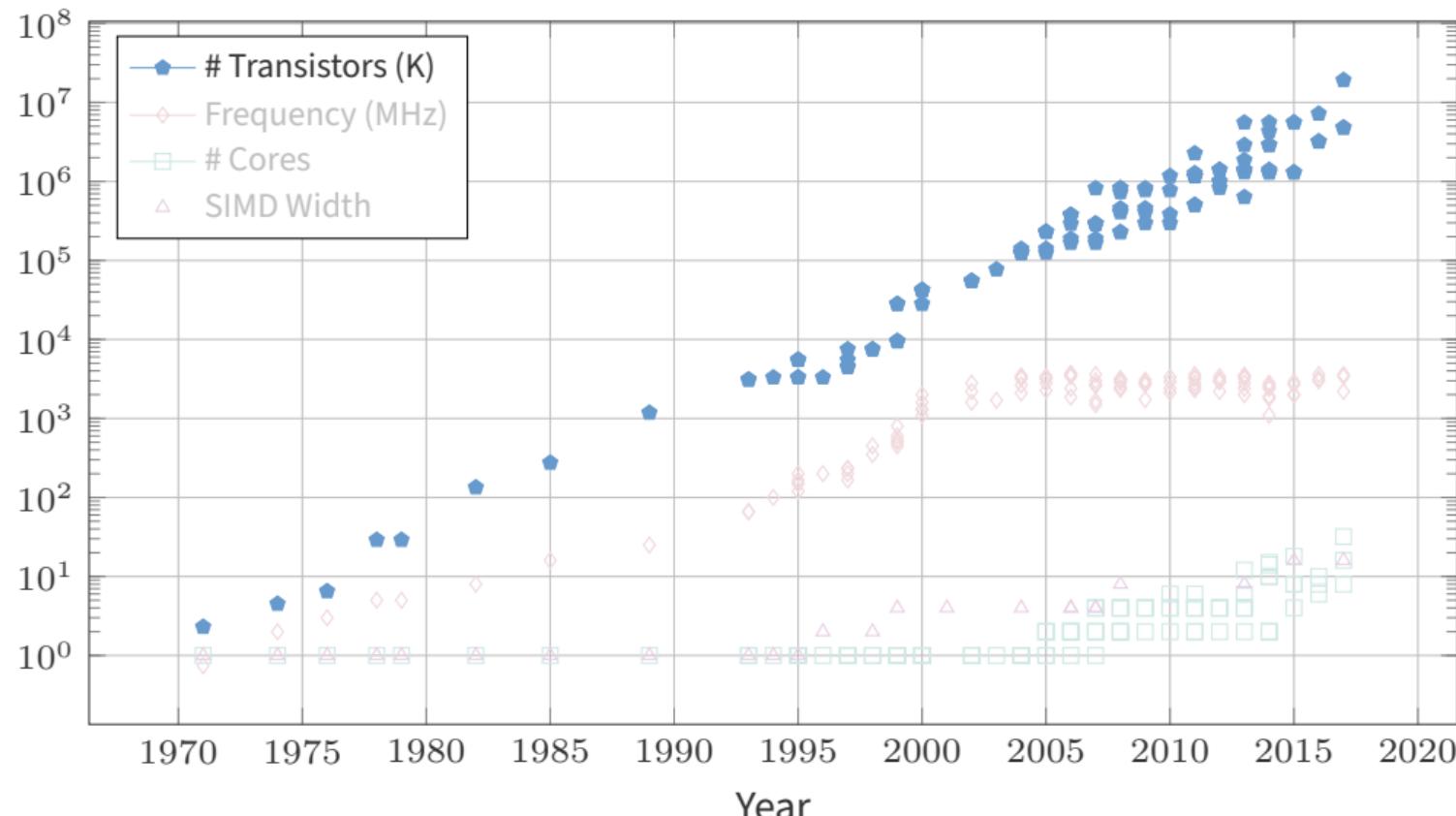
$$24 \text{ h} \cdot \frac{1^{\frac{t}{2}}}{2} = \frac{1}{60} \text{ s}$$

$$t = 2 \cdot \log_2(86400 \cdot 60) \approx 44.6 \text{ years}$$

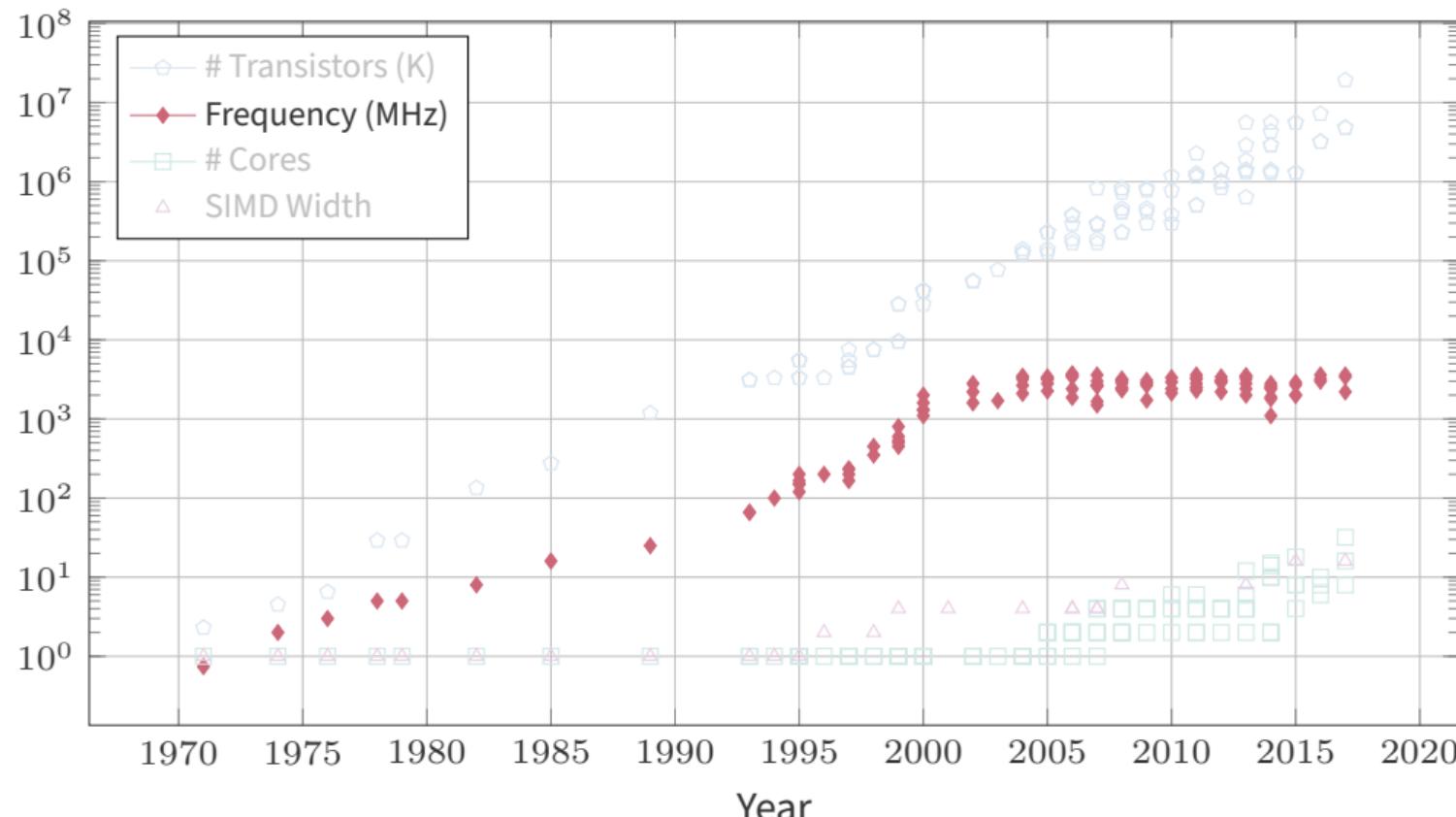
## Moore and more



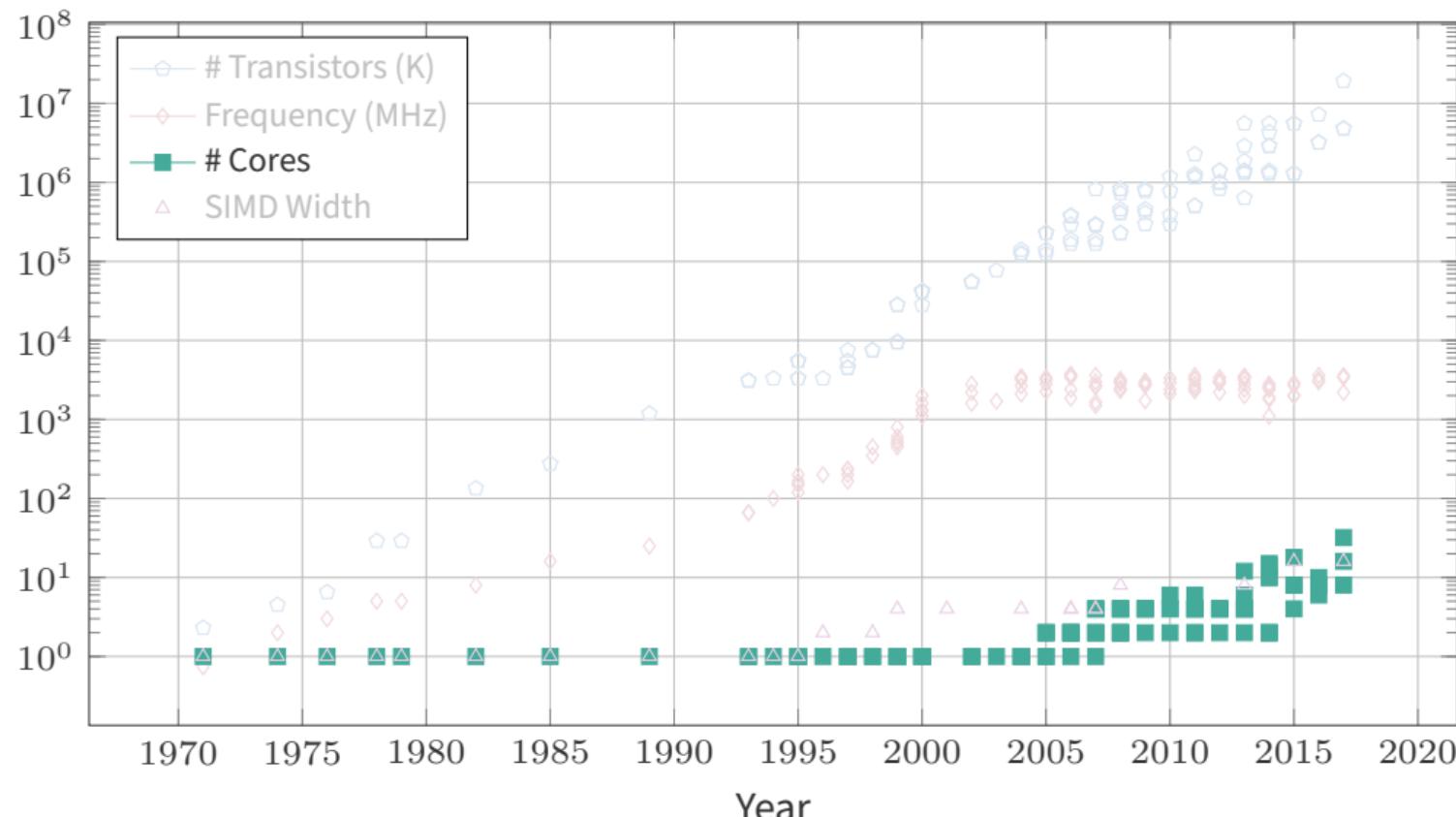
# Moore and more



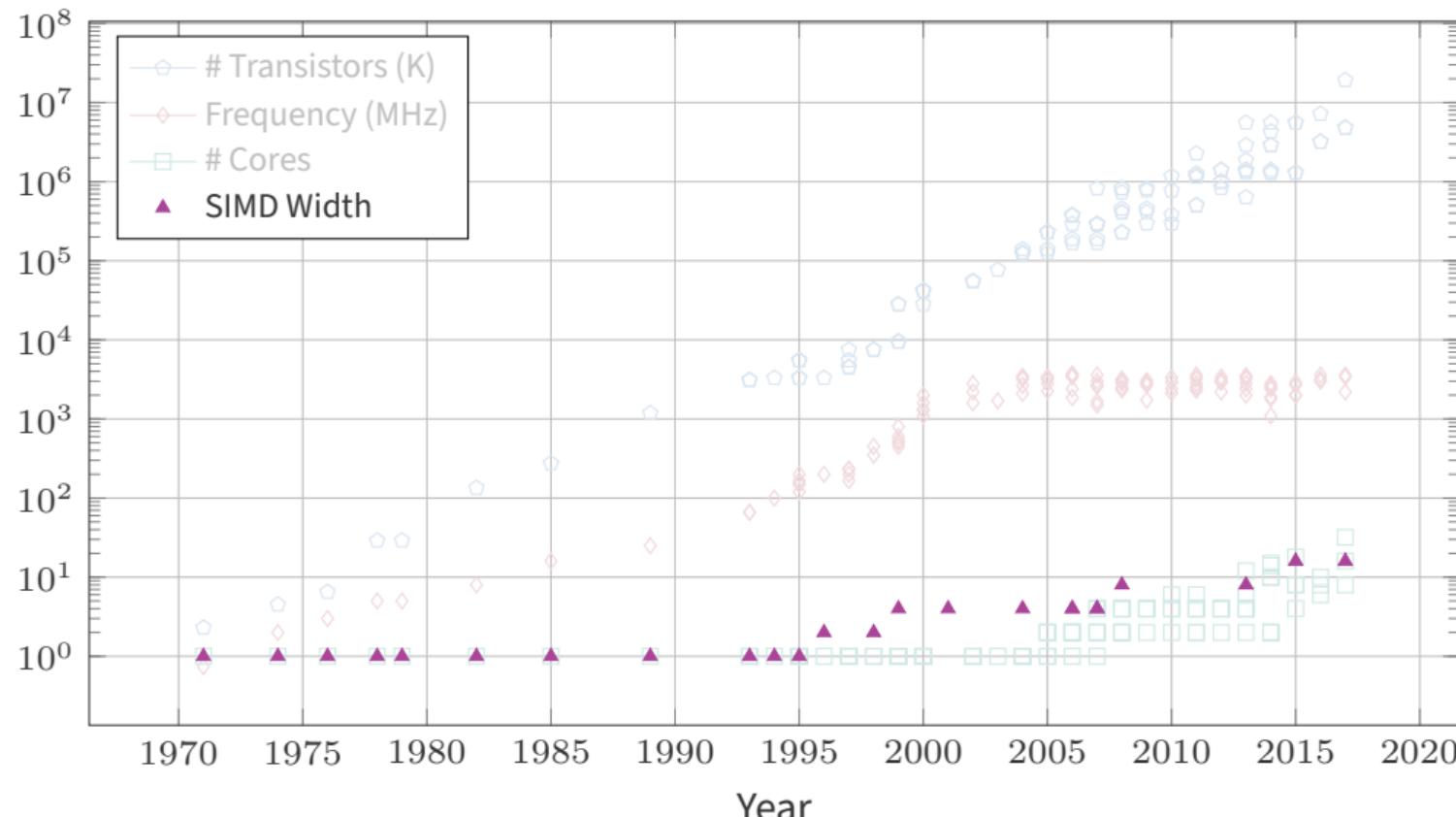
## Moore and more



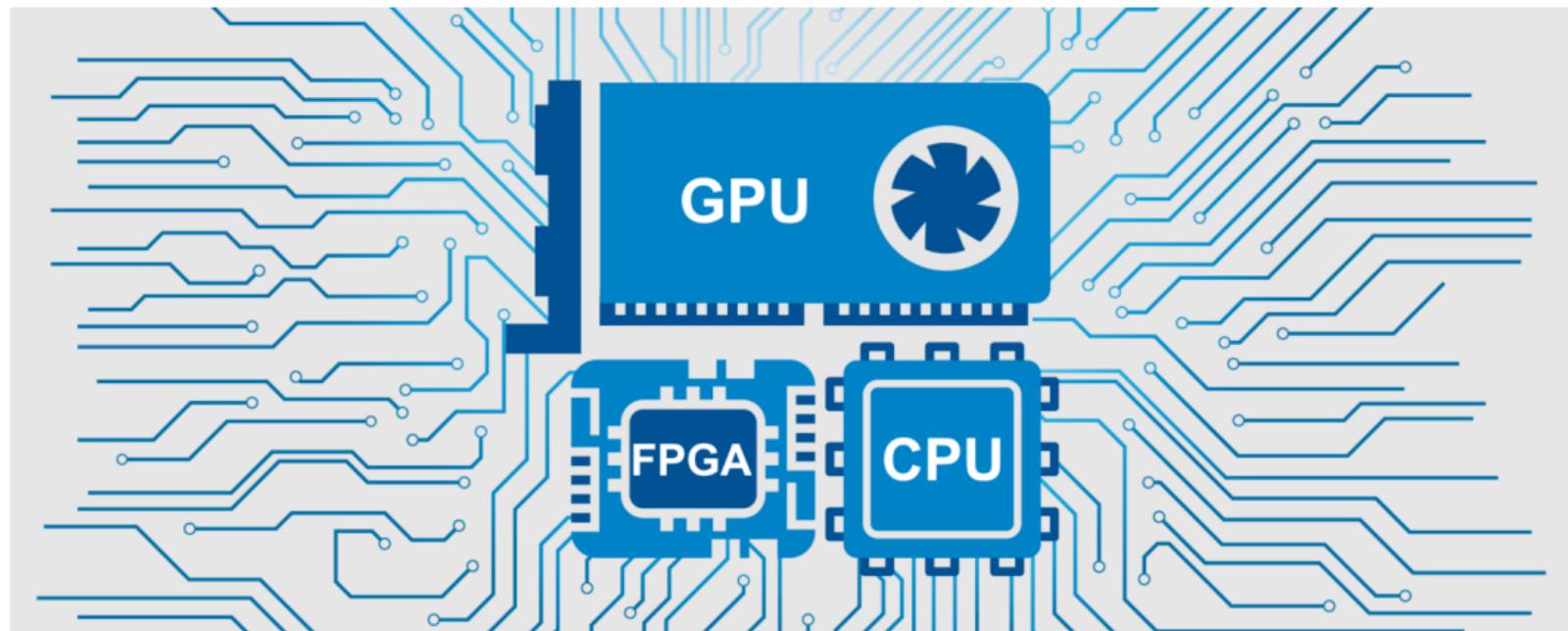
## Moore and more



## Moore and more



# Heterogeneous Hardware



## Training AI

Model	Training end	Chip type	TFLOP/s (max)	Chip count	Wall clock (days)	Total time (years)	Cost (US\$)	MMLU ▲
<b>GPT-3</b>	Apr/2020	V100	130	10,000	15 days	405 years	\$9M	43.9
<b>Llama 1</b>	Jan/2023	A100	312	2,048	21 days	118 years	\$4M	63.4
<b>Llama 2</b>	Jun/2023	A100	312	2,048	35 days	196 years	\$7M	68.0
<b>GPT-4</b>	Aug/2022	A100	312	25,000	95 days	6,507 years	\$224M	86.4
<b>Gemini</b>	Nov/2023	TPUv4	275	57,000	100 days	15,616 years	\$440M	90.0
<b>GPT-5</b>	Apr/2024	H100	989	50,000	120 days	16,438 years	\$612M	
<b>Llama 3</b>	Apr/2024	H100	989					
<b>Olympus</b>	Aug/2024	H100	989					
<b>Gemini 2</b>	Nov/2024	TPUv5	393					

Table. Google DeepMind Gemini training compute (see working, with sources<sup>8</sup>).

## **Structure of a Compiler**

---

# Architecture

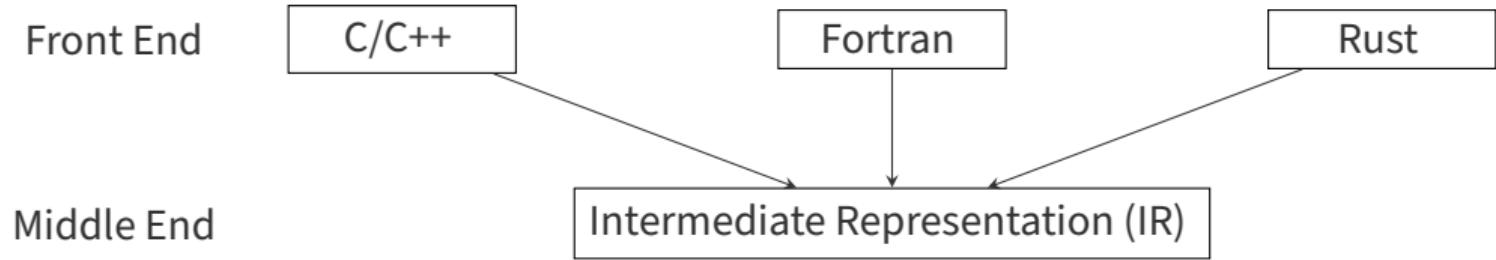
Front End

C/C++

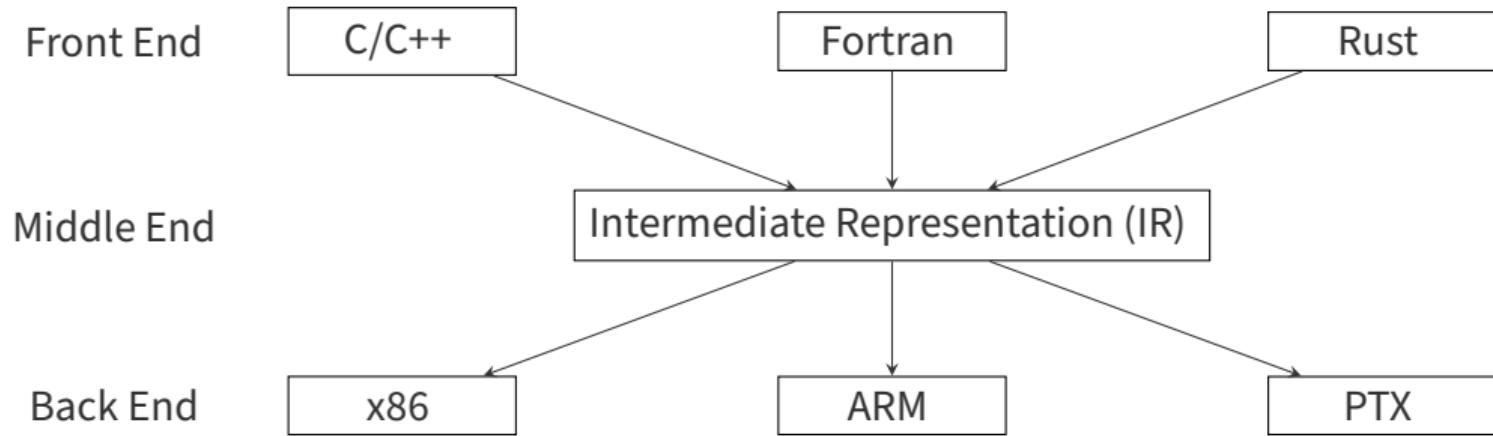
Fortran

Rust

# Architecture



# Architecture



## Running Example

```
float f(float in) {  
    float out;  
    out = 2 * in;  
    return out;  
}
```

## Running Example

```
float f(float in) {  
    float out;  
    out = 2 * in;  
    return out;  
}
```

## Lexical Analysis

- **Input:** Sequence of Characters

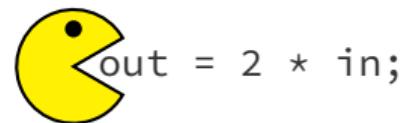
# Lexical Analysis

- **Input:** Sequence of Characters
- **Tasks:**
  - Deconstruct incoming text into **meaningful chunks** called **Tokens**
  - Discard useless characters (whitespace, comments)

# Lexical Analysis

- **Input:** Sequence of Characters
- **Tasks:**
  - Deconstruct incoming text into **meaningful chunks** called **Tokens**
  - Discard useless characters (whitespace, comments)
- **Output:** Sequence of tokens

## Lexical Analysis



## Lexical Analysis

- o  ut = 2 \* in;

## Lexical Analysis

ou  t = 2 \* in;

## Lexical Analysis

$\langle \mathbf{id}, \text{out} \rangle$   = 2 \* in;

## Lexical Analysis

$\langle \mathbf{id}, \text{out} \rangle$   2 \* in;

## Lexical Analysis

$\langle \mathbf{id}, \text{out} \rangle \langle = \rangle$   2 \* in;

## Lexical Analysis

$\langle \mathbf{id}, \text{out} \rangle \langle = \rangle$   \* in;

## Lexical Analysis

$\langle \text{id}, \text{out} \rangle \xrightarrow{=} \langle \text{lit}, 2 \rangle$   \* in;

## Lexical Analysis

$\langle \text{id}, \text{out} \rangle \xrightarrow{=} \langle \text{lit}, 2 \rangle$   in;

## Lexical Analysis

$\langle \text{id}, \text{out} \rangle \langle = \rangle \langle \text{lit}, 2 \rangle \langle * \rangle$   in;

## Lexical Analysis

$\langle \mathbf{id}, \text{out} \rangle \langle = \rangle \langle \mathbf{lit}, 2 \rangle \langle * \rangle \quad \text{in;}$



## Lexical Analysis

$\langle \text{id}, \text{out} \rangle \langle = \rangle \langle \text{lit}, 2 \rangle \langle * \rangle i$   n;

# Lexical Analysis

$\langle \text{id}, \text{out} \rangle \langle = \rangle \langle \text{lit}, 2 \rangle \langle * \rangle \langle \text{id}, \text{in} \rangle$



## Lexical Analysis

`<id, out> <= > <lit, 2> <*> <id, in> <;>`



## Syntax Analysis

- **Input:** Sequence of Tokens

# Syntax Analysis

- **Input:** Sequence of Tokens
- **Tasks**
  - Checks if input adheres to context-free grammar of source language
  - Detects hierarchical structure of input

# Syntax Analysis

- **Input:** Sequence of Tokens
- **Tasks**
  - Checks if input adheres to context-free grammar of source language
  - Detects hierarchical structure of input
- **Output:** Abstract syntax tree (AST)

## Syntax Analysis



$\langle \text{id}, \text{out} \rangle \langle = \rangle \langle \text{lit}, 2 \rangle \langle * \rangle \langle \text{id}, \text{in} \rangle \langle ; \rangle$

# Syntax Analysis



$\langle = \rangle \langle \text{lit}, 2 \rangle \langle * \rangle \langle \text{id}, \text{in} \rangle \langle ; \rangle$

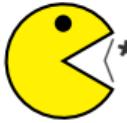
**id** out

## Syntax Analysis

=   $\langle \text{lit}, 2 \rangle \langle * \rangle \langle \text{id}, \text{in} \rangle \langle ; \rangle$

**id** out

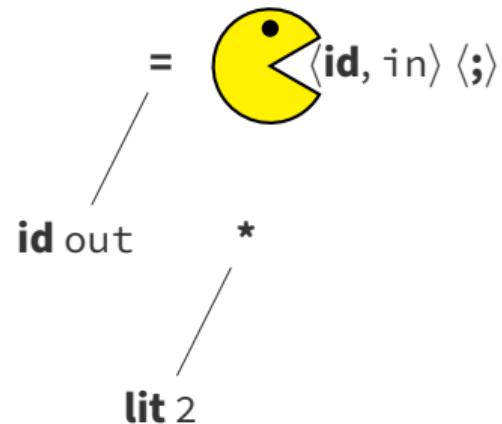
## Syntax Analysis

=   $\langle^* \rangle \langle \mathbf{id}, \mathbf{in} \rangle \langle ; \rangle$

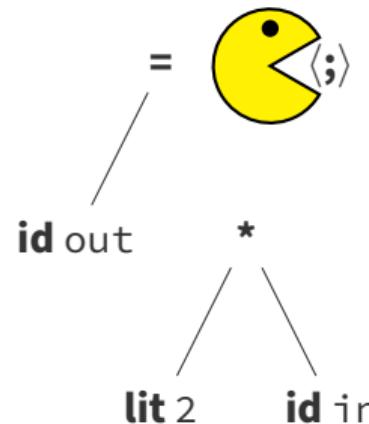
**id** out

**lit** 2

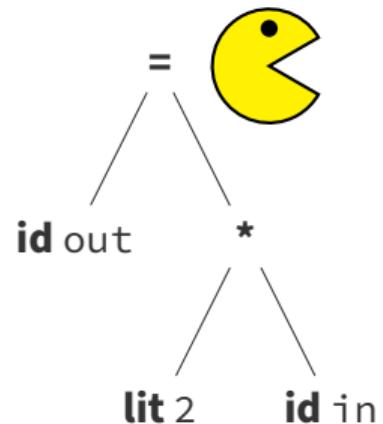
## Syntax Analysis



## Syntax Analysis



## Syntax Analysis



# Semantic Analysis

- **Input:** AST

# Semantic Analysis

- **Input:** AST
- **Tasks:**
  - **Name analysis:** Relate uses to definitions of names
  - **Type analysis:** Elaborate and check the types of entities

# Semantic Analysis

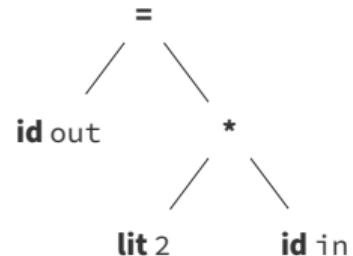
- **Input:** AST
- **Tasks:**
  - **Name analysis:** Relate uses to definitions of names
  - **Type analysis:** Elaborate and check the types of entities
- **Output:** Attributed AST

# Semantic Analysis

```
float f(float in) {  
    float out;  
    out = 2 * in;  
    return out;  
}
```

# Semantic Analysis

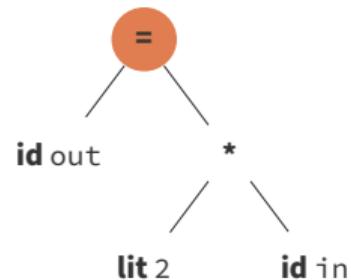
```
float f(float in) {  
    float out;
```



```
    return out;  
}
```

# Semantic Analysis

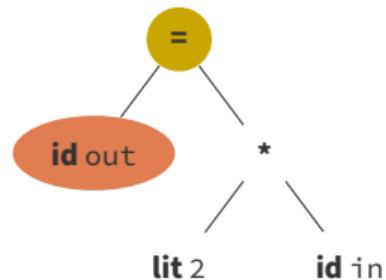
```
float f(float in) {  
    float out;
```



```
    return out;  
}
```

# Semantic Analysis

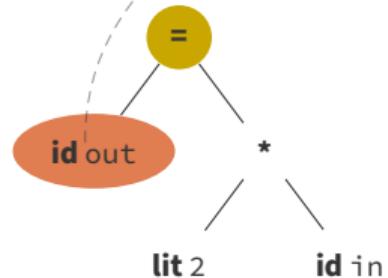
```
float f(float in) {  
    float out;
```



```
return out;  
}
```

# Semantic Analysis

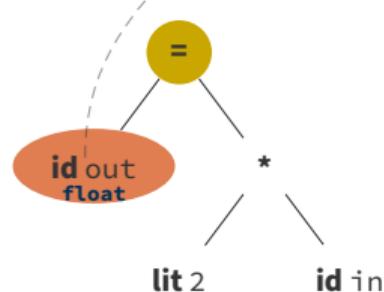
```
float f(float in) {  
    float out;
```



```
    return out;  
}
```

# Semantic Analysis

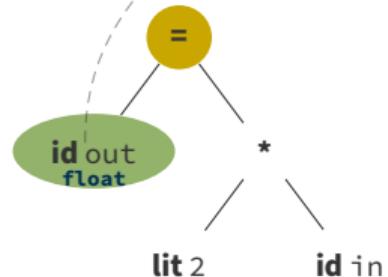
```
float f(float in) {  
    float out;
```



```
    return out;  
}
```

# Semantic Analysis

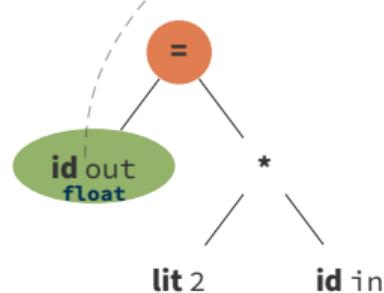
```
float f(float in) {  
    float out;
```



```
    return out;  
}
```

# Semantic Analysis

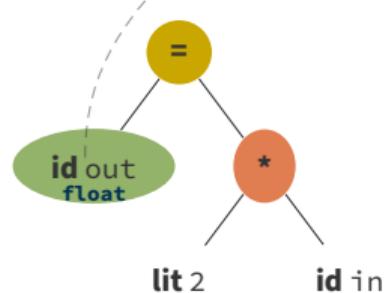
```
float f(float in) {  
    float out;
```



```
    return out;  
}
```

# Semantic Analysis

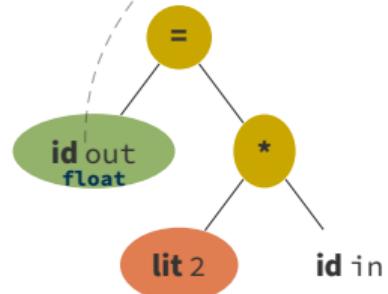
```
float f(float in) {  
    float out;
```



```
    return out;  
}
```

# Semantic Analysis

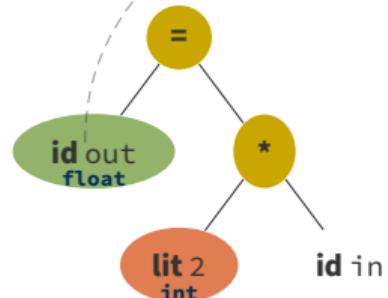
```
float f(float in) {  
    float out;
```



```
    return out;  
}
```

# Semantic Analysis

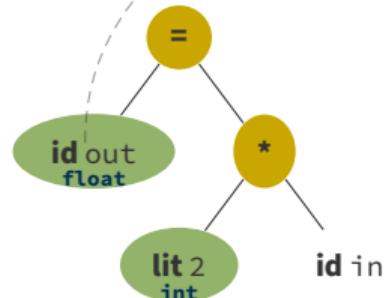
```
float f(float in) {  
    float out;
```



```
        return out;  
}
```

# Semantic Analysis

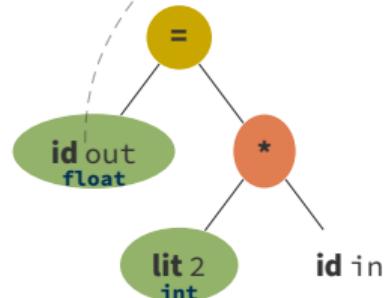
```
float f(float in) {  
    float out;
```



```
    return out;  
}
```

# Semantic Analysis

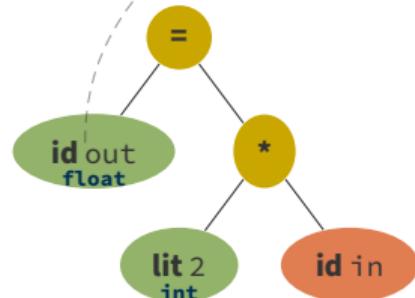
```
float f(float in) {  
    float out;
```



```
    return out;  
}
```

# Semantic Analysis

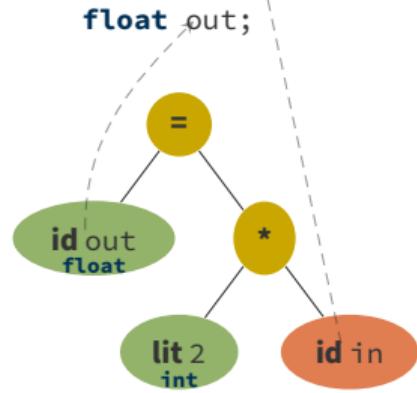
```
float f(float in) {  
    float out;
```



```
    return out;  
}
```

# Semantic Analysis

```
float f(float in) {
```

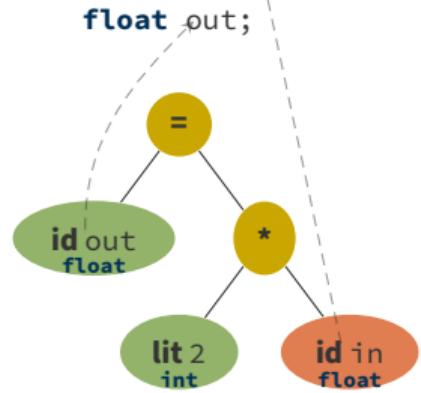


```
    return out;
```

```
}
```

# Semantic Analysis

```
float f(float in) {
```



```
    return out;
```

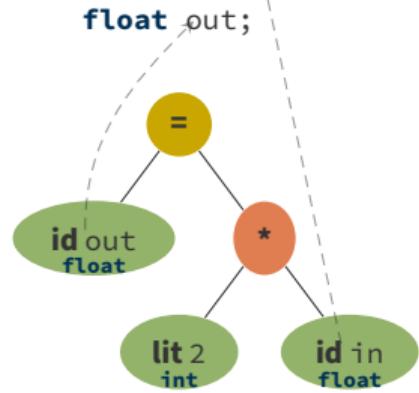
```
}
```

# Semantic Analysis

```
float f(float in) {  
    float out;  
  
    =  
    / \  
   id out float  
      *  
     / \  
    lit 2 int  
    id in float  
  
    return out;  
}
```

# Semantic Analysis

```
float f(float in) {
```

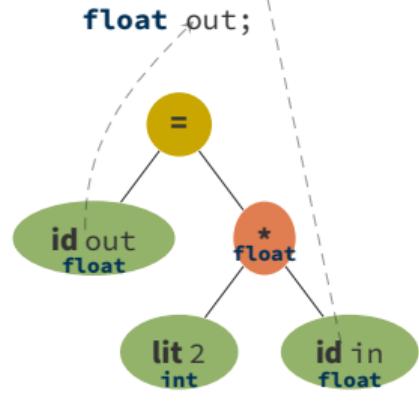


```
    return out;
```

```
}
```

# Semantic Analysis

```
float f(float in) {
```



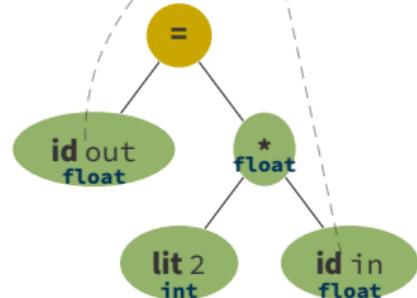
```
    return out;
```

```
}
```

# Semantic Analysis

```
float f(float in) {
```

```
    float out;
```

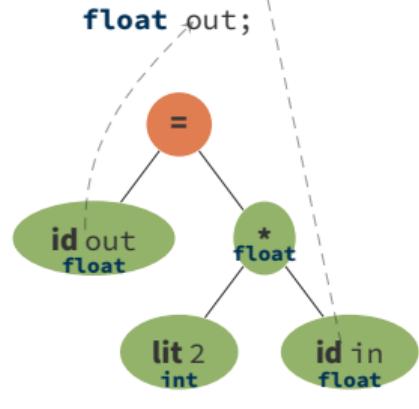


```
    return out;
```

```
}
```

# Semantic Analysis

```
float f(float in) {
```



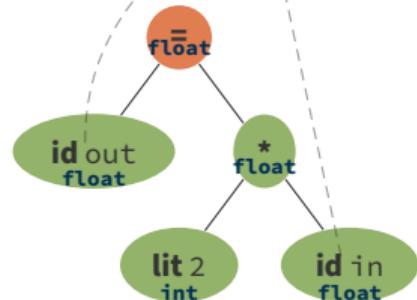
```
    return out;
```

```
}
```

# Semantic Analysis

```
float f(float in) {
```

```
    float out;
```



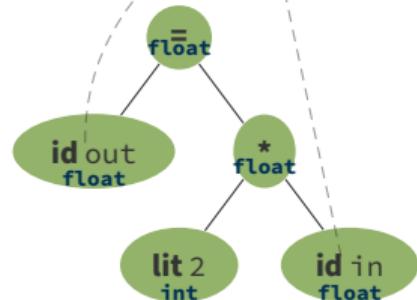
```
    return out;
```

```
}
```

# Semantic Analysis

```
float f(float in) {
```

```
    float out;
```



```
    return out;
```

```
}
```

# Transformation

- **Input:** Attributed AST

## Transformation

- **Input:** Attributed AST
- **Tasks:** Lower high-level language constructs to target code

# Transformation

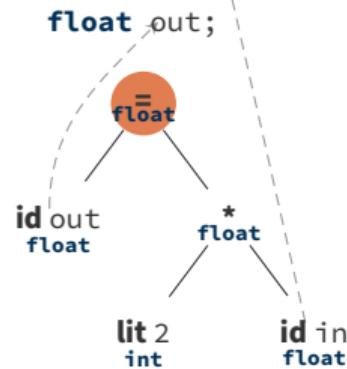
- **Input:** Attributed AST
- **Tasks:** Lower high-level language constructs to target code
- **Output:** Machine code or **intermediate representation (IR)**

# Transformation

```
float f(float in) {  
    float out;  
    float =  
        id out  
        float  
        *  
        float  
        lit 2  
        int  
        id in  
        float  
  
    return out;  
}
```

# Transformation

```
float f(float in) {
```

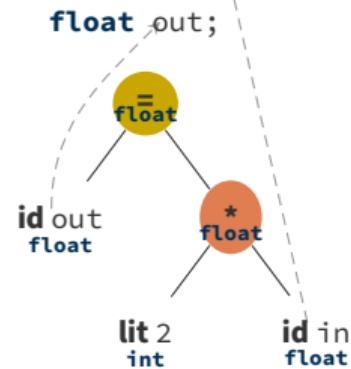


```
    return out;
```

```
}
```

# Transformation

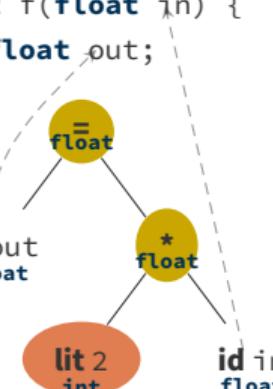
```
float f(float in) {
```



```
    return out;
```

```
}
```

# Transformation

```
float f(float in) {  
    float out;  


```
graph TD; A((float)) -- "*" --> B((float)); A --- C[id out float]; A --- D[lit 2 int];
```



The diagram illustrates an abstract syntax tree (AST) for the expression lit 2 * id in. The root node is a yellow circle labeled float. It has two children: a yellow circle labeled * and another yellow circle labeled float. The left child * has one child, a red oval labeled lit 2 int. The right child float has one child, a blue oval labeled id in float. Dashed lines connect the root node to its children, and dashed lines also connect the children to their respective children. The code above the tree defines a function f that takes a float parameter in and returns a float value out.

  
        return out;  
    }
```

# Transformation

```
float f(float in) {  
    float out;  
  
    id out float  
    * float  
    lit 2 int  
    id in float  
  
    return out;  
}
```

2

# Transformation

```
float f(float in) {  
    float out;  
  
    id out  
    float  
  
    * float  
  
    lit 2  
    int  
  
    return out;  
}
```

2

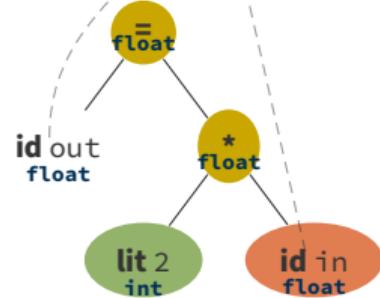
# Transformation

```
float f(float in) {  
    float out;  
  
    id out  
    float  
  
    * float  
  
    lit 2  
    int  
  
    id in  
    float  
  
    return out;  
}
```

t1 = inttofloat 2

# Transformation

```
float f(float in) {  
    float out;
```



```
        return out;
```

```
}
```

```
t1 = inttofloat 2
```

# Transformation

```
float f(float in) {  
    float out;  
  
    id out  
    float  
  
    * float  
  
    lit 2 int  
    id in float  
  
    return out;  
}
```

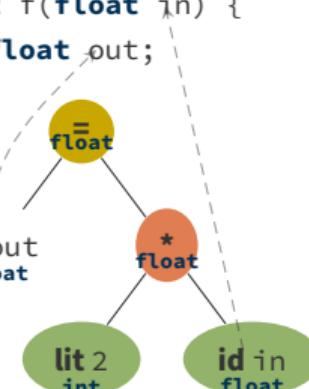
```
t1 = inttofloat 2  
t2 = load in
```

# Transformation

```
float f(float in) {  
    float out;  
  
    id out  
    float  
  
    * float  
  
    lit 2 int  
    id in float  
  
    return out;  
}
```

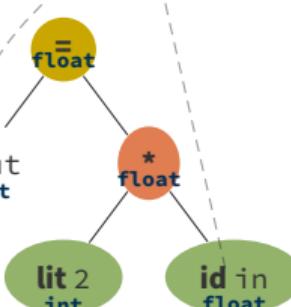
```
t1 = inttofloat 2  
t2 = load in
```

# Transformation

```
float f(float in) {  
    float out;  
  
      
    *  
    / \   
   id in float  lit 2 int  
  
    return out;  
}
```

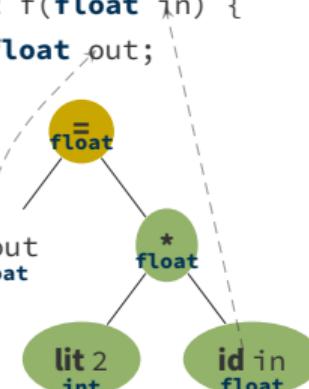
```
t1 = inttofloat 2  
t2 = load in
```

# Transformation

```
float f(float in) {  
    float out;  
  
      
    return out;  
}
```

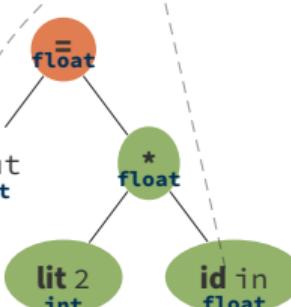
```
t1 = inttofloat 2  
t2 = load in  
t3 = mul t1, t2
```

# Transformation

```
float f(float in) {  
    float out;  
  
      
    return out;  
}
```

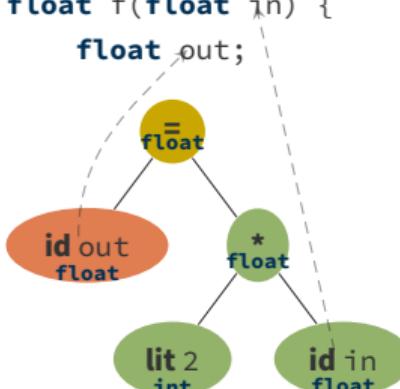
```
t1 = inttofloat 2  
t2 = load in  
t3 = mul t1, t2
```

# Transformation

```
float f(float in) {  
    float out;  
  
      
    return out;  
}
```

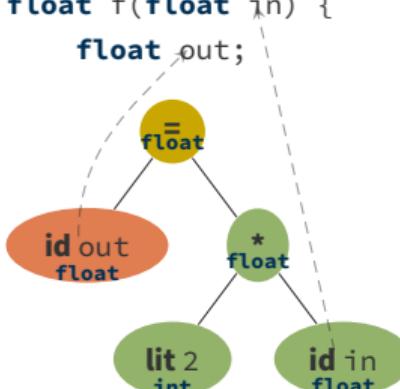
```
t1 = inttofloat 2  
t2 = load in  
t3 = mul t1, t2
```

# Transformation

```
float f(float in) {  
    float out;  
  
      
    return out;  
}
```

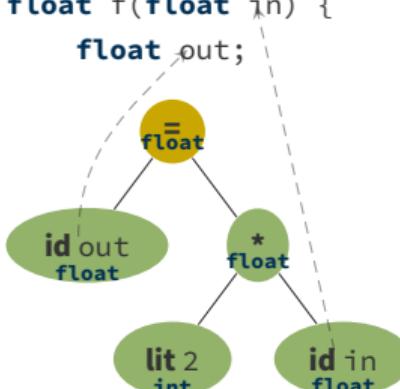
```
t1 = inttofloat 2  
t2 = load in  
t3 = mul t1, t2
```

# Transformation

```
float f(float in) {  
    float out;  
  
      
    return out;  
}
```

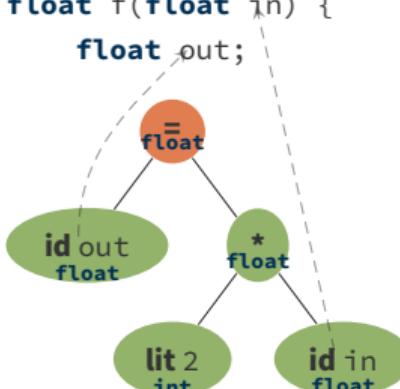
```
t1 = inttofloat 2  
t2 = load in  
t3 = mul t1, t2  
      out
```

# Transformation

```
float f(float in) {  
    float out;  
  
      
    return out;  
}
```

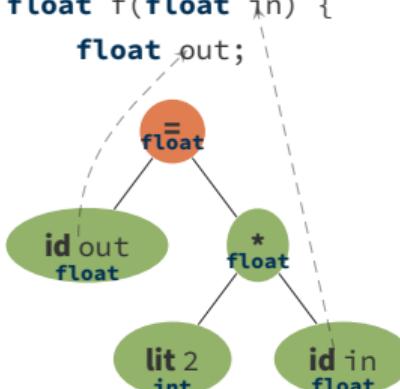
```
t1 = inttofloat 2  
t2 = load in  
t3 = mul t1, t2  
          out
```

# Transformation

```
float f(float in) {  
    float out;  
  
      
    return out;  
}
```

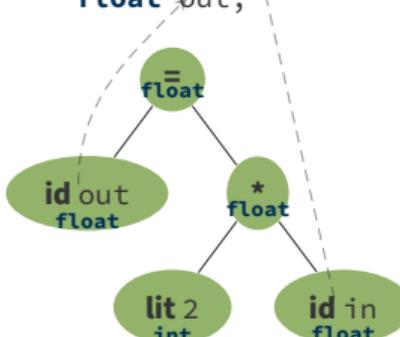
```
t1 = inttofloat 2  
t2 = load in  
t3 = mul t1, t2  
          out
```

# Transformation

```
float f(float in) {  
    float out;  
  
      
    return out;  
}
```

```
t1 = inttofloat 2  
t2 = load in  
t3 = mul t1, t2  
store t3, out
```

# Transformation

```
float f(float in) {  
    float out;  
  
      
    return out;  
}
```

```
t1 = inttofloat 2  
t2 = load in  
t3 = mul t1, t2  
store t3, out
```

# Optimization

- **Input:** IR

# Optimization

- **Input:** IR
- **Tasks:**
  - Reduce resource consumption: execution time, memory
  - Preserve semantics

# Optimization

- **Input:** IR
- **Tasks:**
  - Reduce resource consumption: execution time, memory
  - Preserve semantics
- **Output:** IR

# Optimization

```
t1 = inttofloat 2
t2 = load in
t3 = mul t1, t2
store t3, out
```

# Optimization

```
t1 = 2.f
t2 = load in
t3 = mul t1, t2
store t3, out
```

# Optimization

```
t2 = load in  
t3 = mul 2.f, t2  
store t3, out
```

# Optimization

```
t2 = load in  
t3 = add t2, t2  
store t3, out
```

# Code Generation

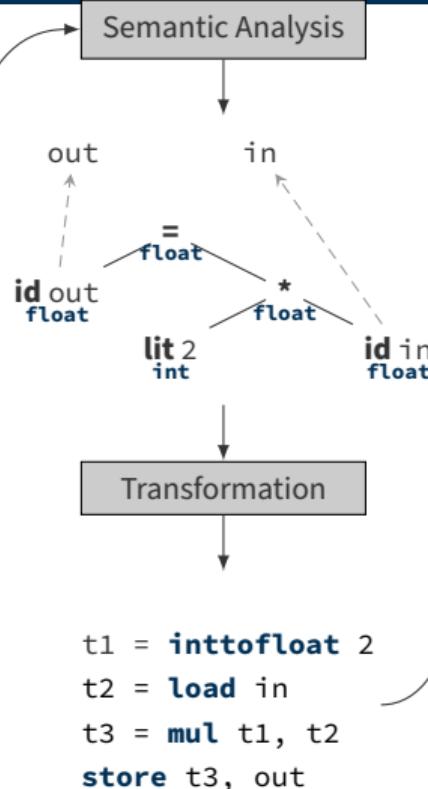
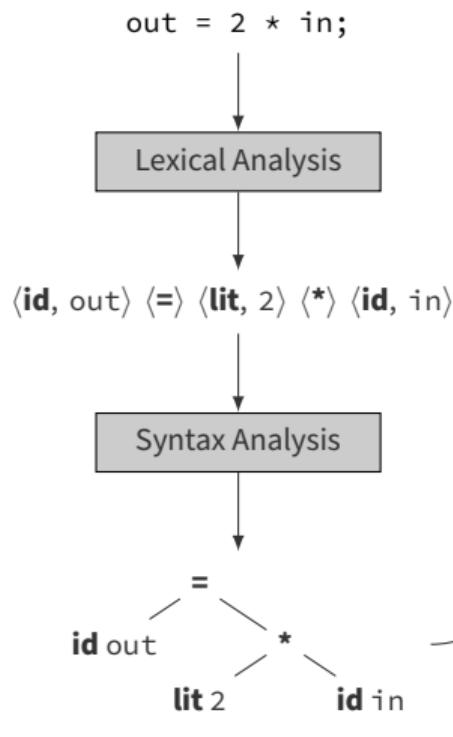
- **Input:** IR
- **Tasks:**
  - Instruction selection: use ISA of target arch
  - Scheduling: re-order instructions
  - Register allocation: allocate registers to IR vars
- **Output:** Machine code

# Code Generation

```
t2 = load in  
t3 = add t2, t2  
store t3, out
```

```
movss in, %xmm0  
addss %xmm0, %xmm0  
movss %xmm0, out
```

# Summary



t2 = **load** in  
t3 = **add** t2, t2  
**store** t3, out

Code Generation

**movss** in, %xmm0  
**adss** %xmm0, %xmm0  
**movss** %xmm0, out

## **Domain-Specific Languages**

---

# General Purpose vs. Domain-Specific Languages

## General Purpose Language

- Java
- C/C++
- Rust
- Fortran
- Python
- Perl

## Domain-Specific Languages

- HTML
- SQL
- MATLAB
- RenderMan, GLSLang, HLSL, ...
- GraphViz
- MLL

# Embedded DSLs

## SQL

```
SELECT name, age  
FROM person  
WHERE age >= 18
```

## C++ → SQL

```
string sql =  
"SELECT name, age  
      FROM person  
     WHERE age >= 18";  
mysql_query(conn, sql);
```

## C++ → sqlpp11 (Embedded DSL)

```
select(person.name, person.age)  
      .from(person)  
     .where(person.age >= 18)
```

## Embedded DSLs

- SQL libraries
- Halide
- PyTorch
- TensorFlow
- JAX
- Z3 (SMT solver)
- MLL

## Not DSL

```
class Vec2:  
    def __init__(self, x, y):  
        self.x = x  
        self.y = y  
  
    def __add__(self, v):  
        return Vec2(self.x + v.x, self.y + v.y)
```

```
r = v + w
```

# Embedded DSL

```
class Vec2:  
    def __init__(self, x, y):  
        self.x = x  
        self.y = y  
  
    def __add__(self, v):  
        return Add(self, v)  
  
class Add:  
    def __init__(self, lhs, rhs):  
        self.lhs = lhs  
        self.rhs = rhs  
  
r = v + w
```

### Scaled Exponential Linear Units (SELU)

```
def selu(x, alpha=1.67, lmbda=1.05):
    return lmbda * jnp.where(x > 0, x, alpha * jnp.exp(x) - alpha)

x = random.normal(key, (1000000,))
%timeit selu(x).block_until_ready()
```

---

```
100 loops, best of 5: 2.05 ms per loop
```

## Scaled Exponential Linear Units (SELU)

```
def selu(x, alpha=1.67, lmbda=1.05):
    return lmbda * jnp.where(x > 0, x, alpha * jnp.exp(x) - alpha)

x = random.normal(key, (1000000,))
%timeit selu(x).block_until_ready()
```

---

100 loops, best of 5: 2.05 ms per loop

## Just-In-Time Compilation (JIT)

```
selu_jit = jax.jit(selu)
selu_jit(x).block_until_ready() # Warm up

%timeit selu_jit(x).block_until_ready()
```

---

10000 loops, best of 5: 150 µs per loop

## Grad

```
def sum_logistic(x):
    return jnp.sum(1.0 / (1.0 + jnp.exp(-x)))

x_small = jnp.arange(3.)
derivative_fn = grad(sum_logistic)
print(derivative_fn(x_small))
```

---

```
[0.25      0.19661194  0.10499357]
```

## Grad

```
def sum_logistic(x):
    return jnp.sum(1.0 / (1.0 + jnp.exp(-x)))

x_small = jnp.arange(3.)
derivative_fn = grad(sum_logistic)
print(derivative_fn(x_small))
```

---

```
[0.25      0.19661194  0.10499357]
```

## Grad + JIT

```
derivative_fn_jit = jax.jit(grad(sum_logistic))
```

## Logistic Function

$$\frac{1}{1 + \exp(-x)}$$

## Java

```
var dag      = new DAG();
var logistic = dag.lit1().div(dag.lit1().add(x().neg().exp()));
var derivative = logistic.backwards();
derivative.llvm("logistic.ll"); // "JIT"
```

## Logistic Function

$$\frac{1}{1 + \exp(-x)}$$

### Java

```
var dag      = new DAG();
var logistic = dag.lit1().div(dag.lit1().add(x().neg().exp()));
var derivative = logistic.backwards();
derivative.llvm("logistic.ll"); // "JIT"
```

### Python

```
dag      = mll.DAG()
logistic = 1.0 / (1.0 + (dag.exp(-dag.x)))
derivative = logistic.backwards()
derivative.llvm("logistic.ll") # "JIT"
```

# **LLVM**

---

- open source
- large, active research community
- used in industry:  
Apple, Google, Intel, NVIDIA, Sony, ...
- front-ends for many languages:  
C/C++, Fortran, Rust, Swift, Julia, Haskell, ...
- back-ends for many architectures:  
X86(-64), ARM, RISC-V, MIPS, WebAssembly, ...
- **HUGE**



- open source
  - large, active research community
  - used in industry:  
Apple, Google, Intel, NVIDIA, Sony, ...
  - front-ends for many languages:  
C/C++, Fortran, Rust, Swift, Julia, Haskell, ...
  - back-ends for many architectures:  
X86(-64), ARM, RISC-V, MIPS, WebAssembly, ...
- . **HUGE**



- open source
- large, active research community
- used in industry:  
Apple, Google, Intel, NVIDIA, Sony, ...
- front-ends for many languages:  
C/C++, Fortran, Rust, Swift, Julia, Haskell, ...
- back-ends for many architectures:  
X86(-64), ARM, RISC-V, MIPS, WebAssembly, ...

. **HUGE**



## LLVM Intermediate Representation

- SSA-based representation of control flow graphs
- dumpable in human-readable, assembly-like form (\*.ll)
- dumpable as compact bitcode (\*.bc)

**MLL: generate human-readable \*.ll file**

# LLVM Intermediate Representation – Example

```
define i32 @fac(i32 %n) {
entry:
  br label %while-header

while-header:
  %it = phi i32 [ %n, %entry ], [ %it_new, %while-body ]
  %res = phi i32 [ 1, %entry ], [ %res_new, %while-body ]
  %while-condition = icmp ne i32 %it, 0
  br i1 %while-condition, label %while-body, label %while-end

while-body:
  %res_new = mul i32 %res, %it
  %it_new = sub i32 %it, 1
  br label %while-header

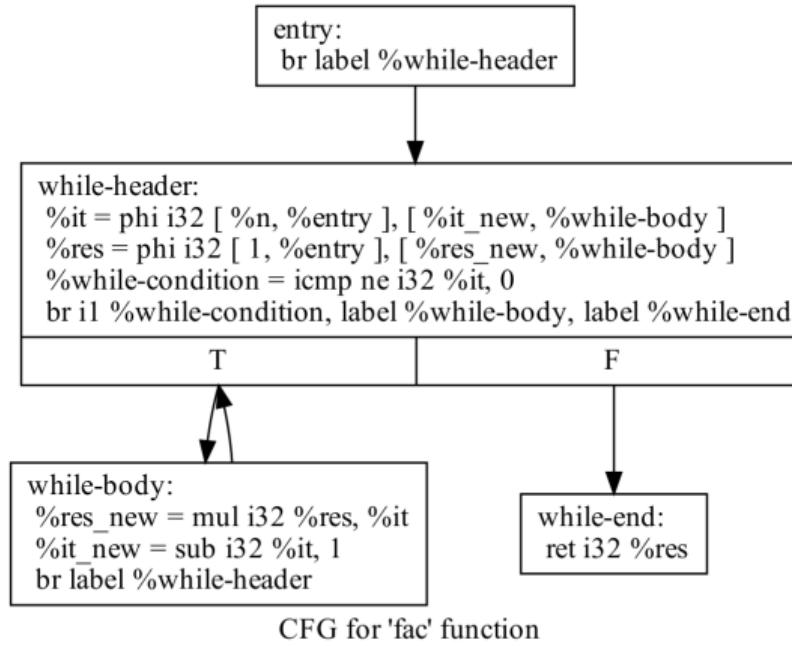
while-end:
  ret i32 %res
}
```

```
int fac(int n) {
  int it = n;
  int res = 1;

  while (it != 0) {
    res = res * it;
    it = it - 1;
  }

  return res;
}
```

# LLVM Intermediate Representation – Example



CFG for 'fac' function

## Useful Commands

- Generate (human readable) LLVM-IR from C/C++ input:

```
clang -Xclang -disable-OO-optnone -emit-llvm -c -S -o OUT.ll IN.c
```

- Draw CFG of function foo from dumped LLVM-IR module:

```
opt -dot-cfg IN.ll; dot -Tpdf cfg.foo.dot > OUT.pdf
```

requires dot/graphviz

- Execute dumped LLVM-IR module:

```
lli IN.ll <argv arguments>
```

- Create binary from dumped LLVM-IR module:

```
clang -o OUT IN.ll
```

- Create architecture-specific assembly:

```
llc -o OUT.s IN.ll
```

- Create binary from architecture specific-assembly:

```
cc -o OUT IN.s
```

- Get more help:

```
<TOOL> --help
```

# Getting Help

- General language reference manual:

<http://llvm.org/docs/LangRef.html>

- Doxygen code documentation:

(well accessible via Google/Bing/DuckDuckGo/...)

<http://llvm.org/doxygen/index.html>

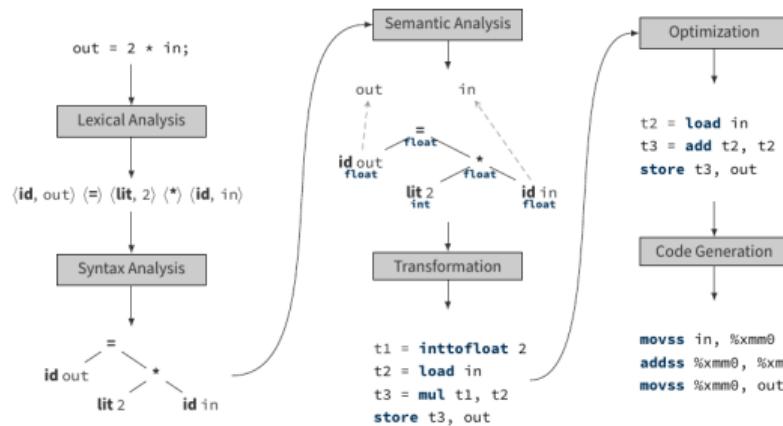
- Full command line tools guide:

<http://llvm.org/docs/CommandGuide/>

## **Conclusion**

---

# Conclusion



SQL

```
SELECT name, age  
FROM person  
WHERE age >= 18;
```

C++ → SQL

```
string sql =  
"SELECT name, age  
FROM person  
WHERE age >= 18";  
mysql_query(conn, sql);
```

C++ → sqlpp11 (Embedded DSL)

```
select(person.name, person.age)  
.from(person)  
.where(person.age >= 18)
```

Thank You!  
Questions?