

regex

Matt Might
University of Utah
matt.might.net



source: xkcd.com

Join the google group!

Extra TA: Petey Aldous

Today

- Research advertisement
- What is lexical analysis?
- What's a regular expression?

Help wanted



Problems

WANTED

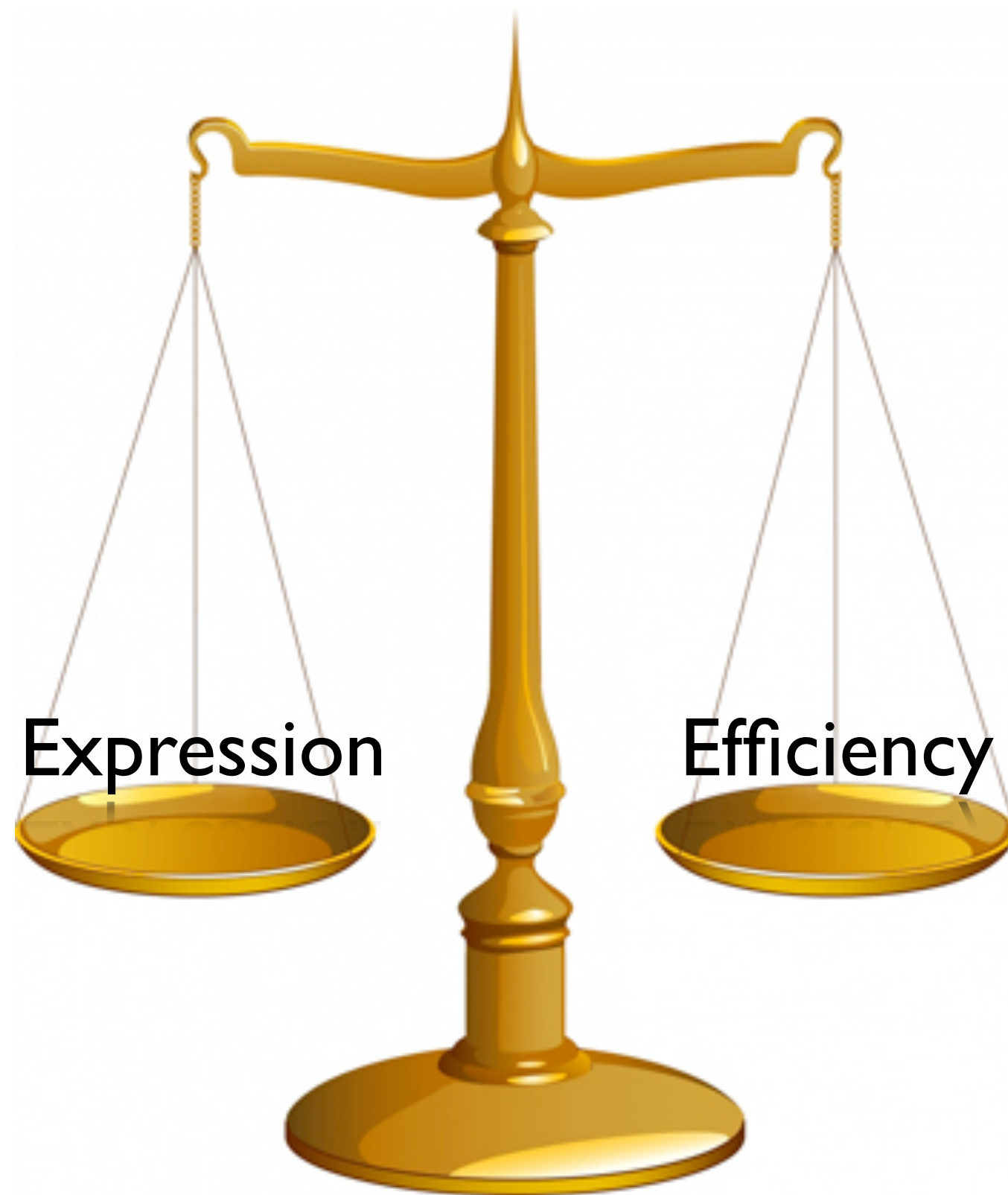
A language that accepts:

$$\begin{aligned}\frac{\partial \rho}{\partial t} &= -\nabla \cdot (\rho \mathbf{u}), \\ \frac{\partial(\rho \mathbf{u})}{\partial t} &= -\nabla \cdot (\rho \mathbf{u} \otimes \mathbf{u}) + \nabla \cdot \boldsymbol{\tau} - \nabla p + \rho \mathbf{g}, \\ \frac{\partial(\rho e_0)}{\partial t} &= -\nabla \cdot [\mathbf{u}(\rho e_0 + p)] + \nabla \cdot (\boldsymbol{\tau} \cdot \mathbf{u}) - \nabla \cdot \mathbf{q} + \rho \mathbf{g} \cdot \mathbf{u}, \\ \frac{\partial(\rho Y_i)}{\partial t} &= -\nabla \cdot (\rho Y_i \mathbf{u}) - \nabla \cdot \mathbf{J}_i + W_i \omega_i,\end{aligned}$$

but performs like this:

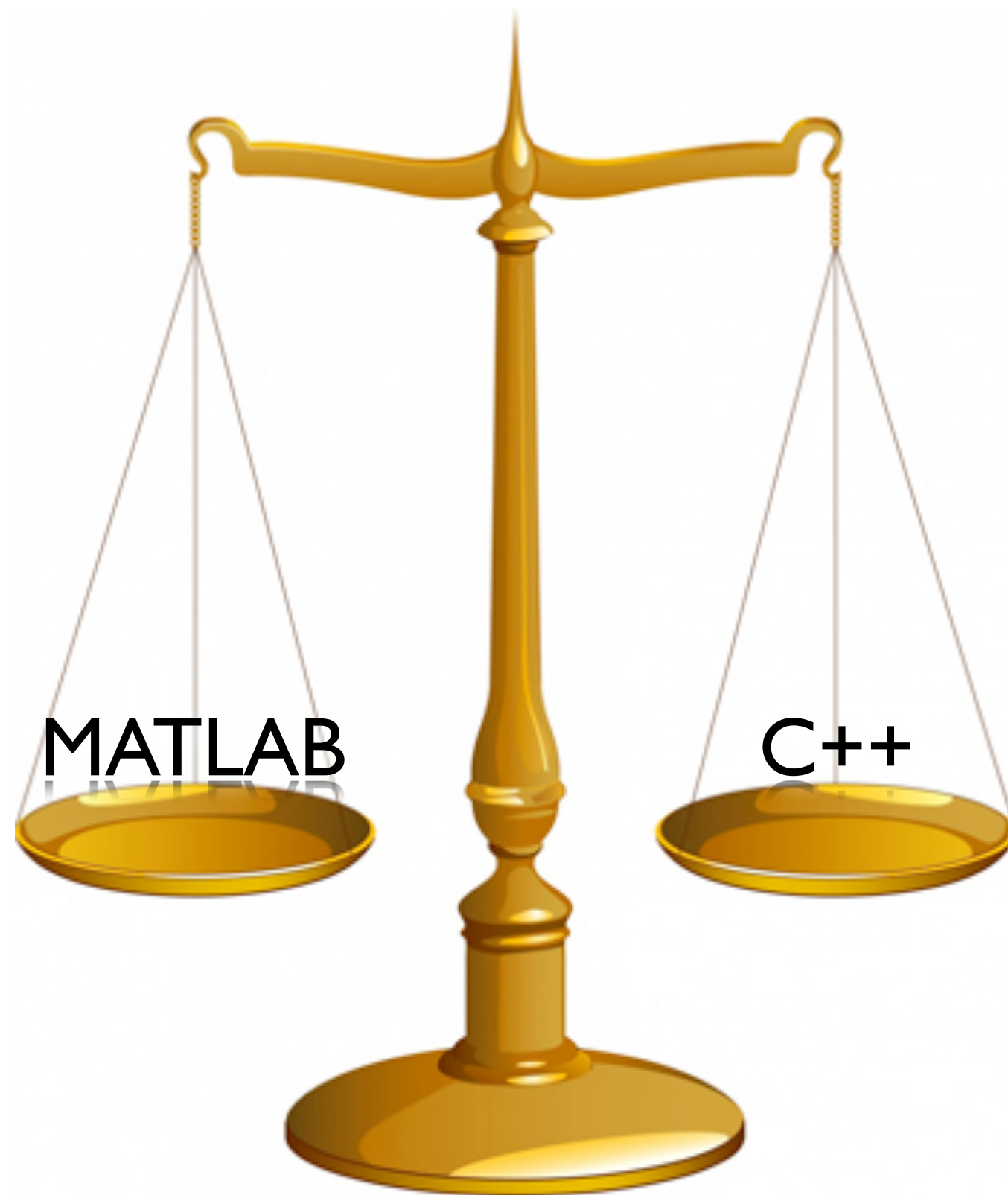
```
Field::const iterator ia = a.begin();
Field::const iterator ib = b.begin();
for(Field::iterator ic = c.begin();
    ic != c.end(); ++ic, ++ia, ++ib) {
    *ic = *ia + sin(*ib);
};
```

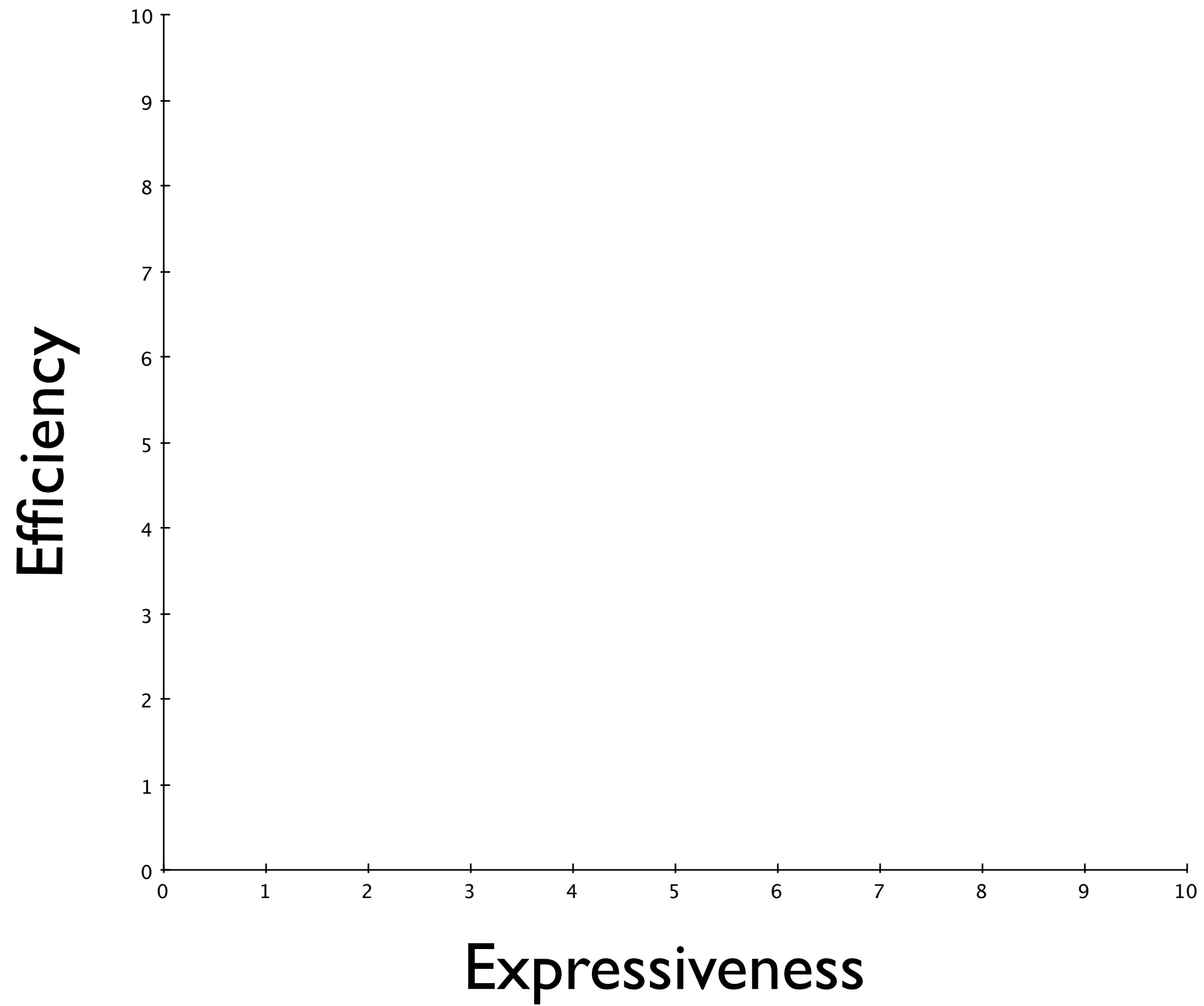


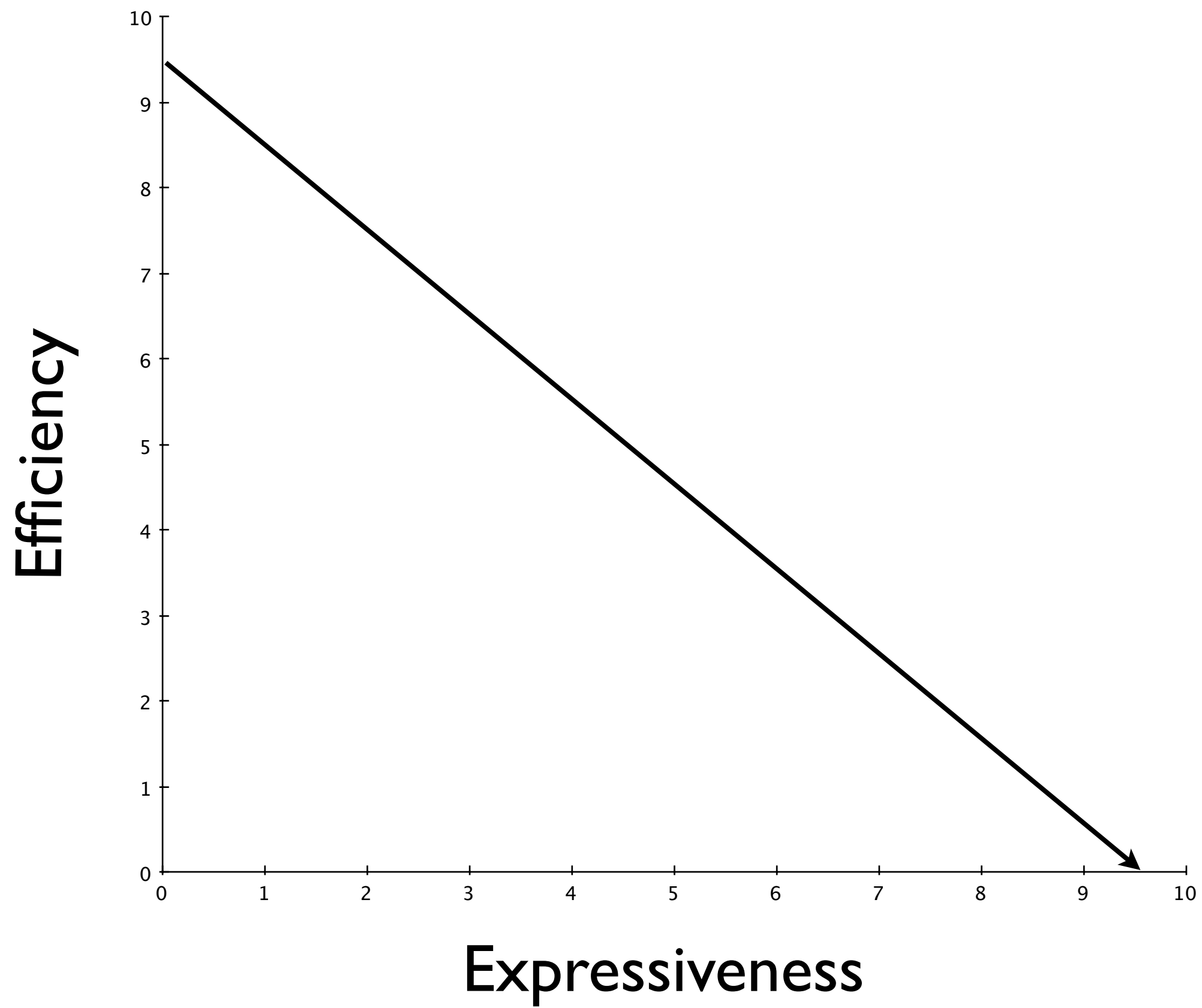



Expression

Efficiency

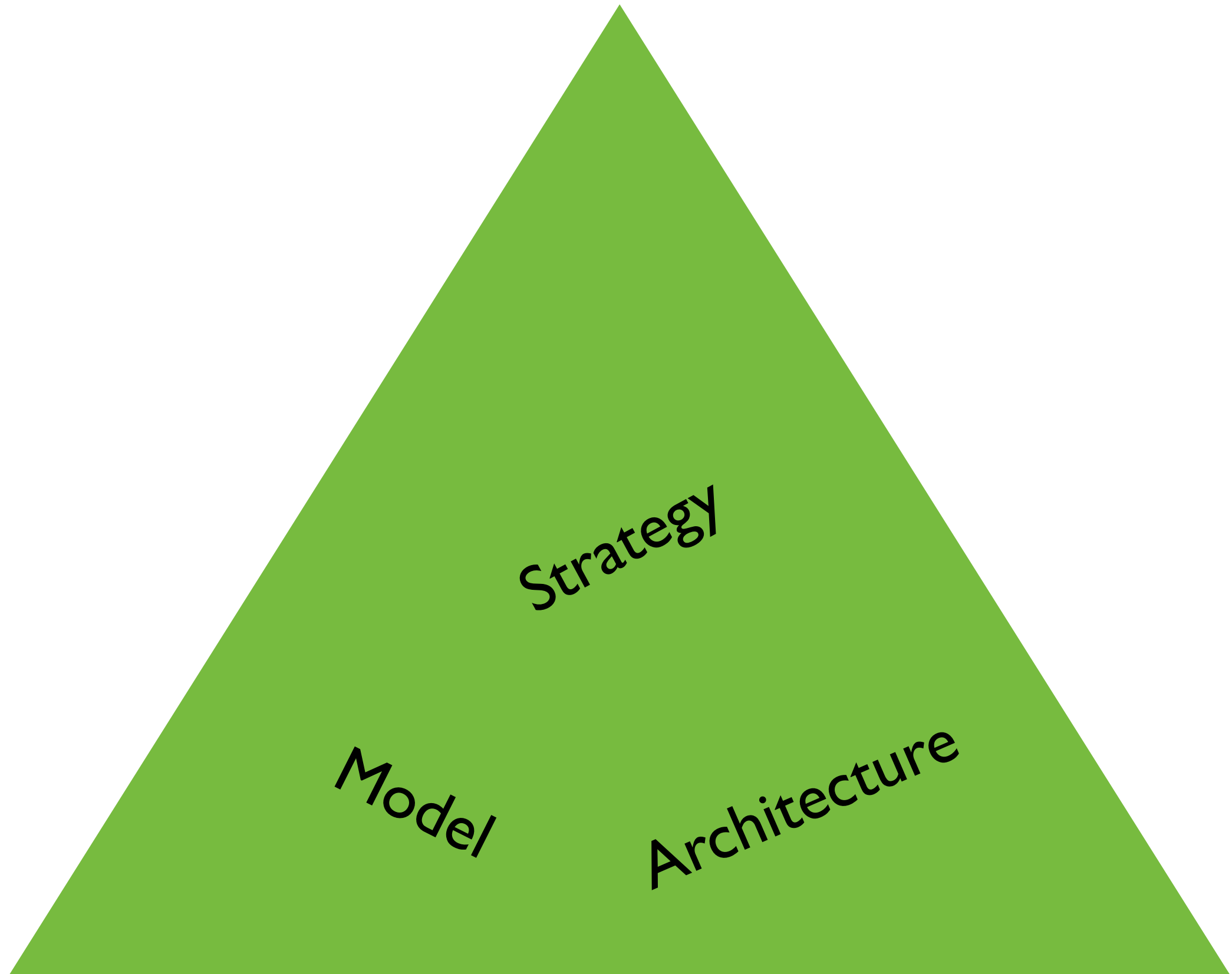




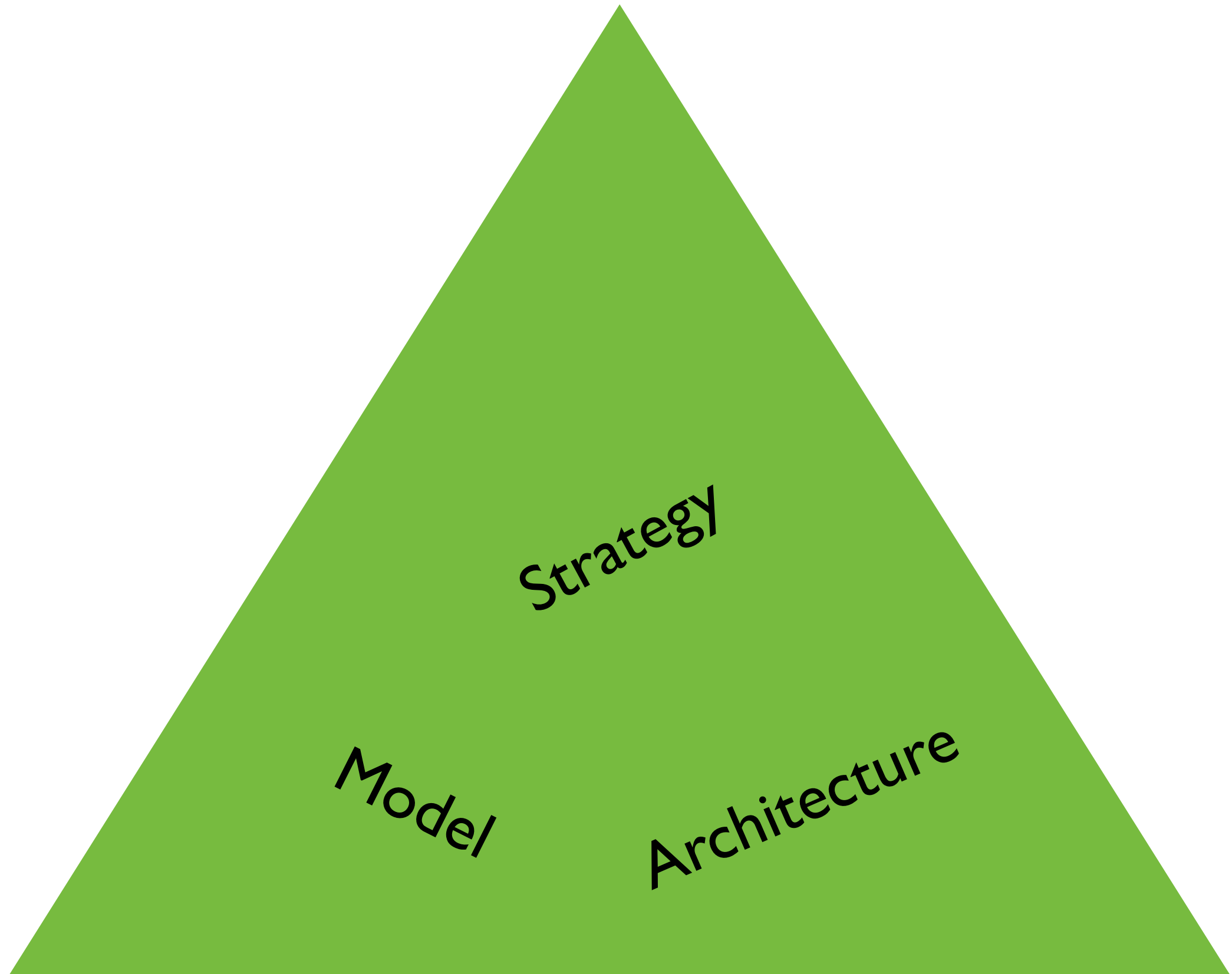


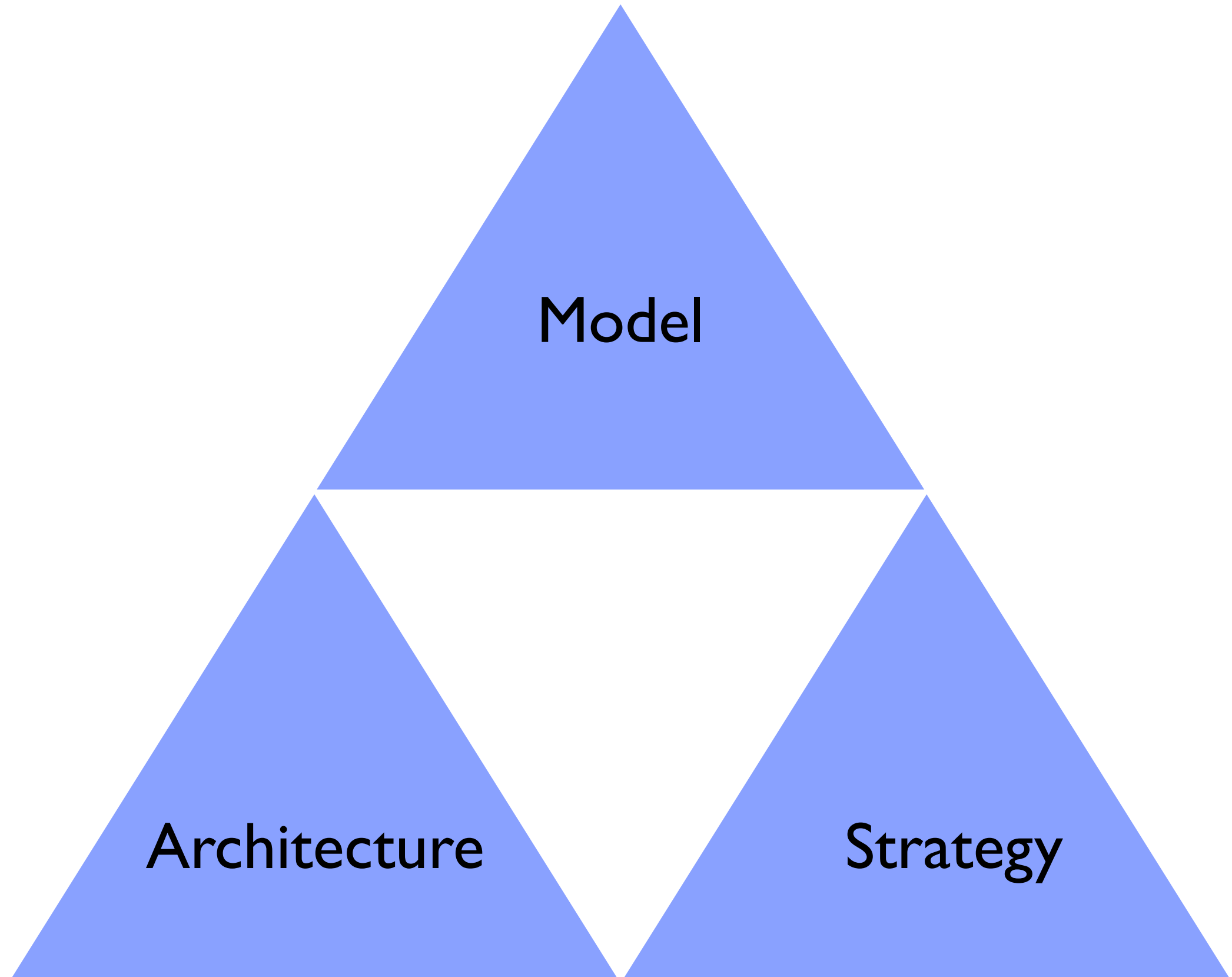
Also...

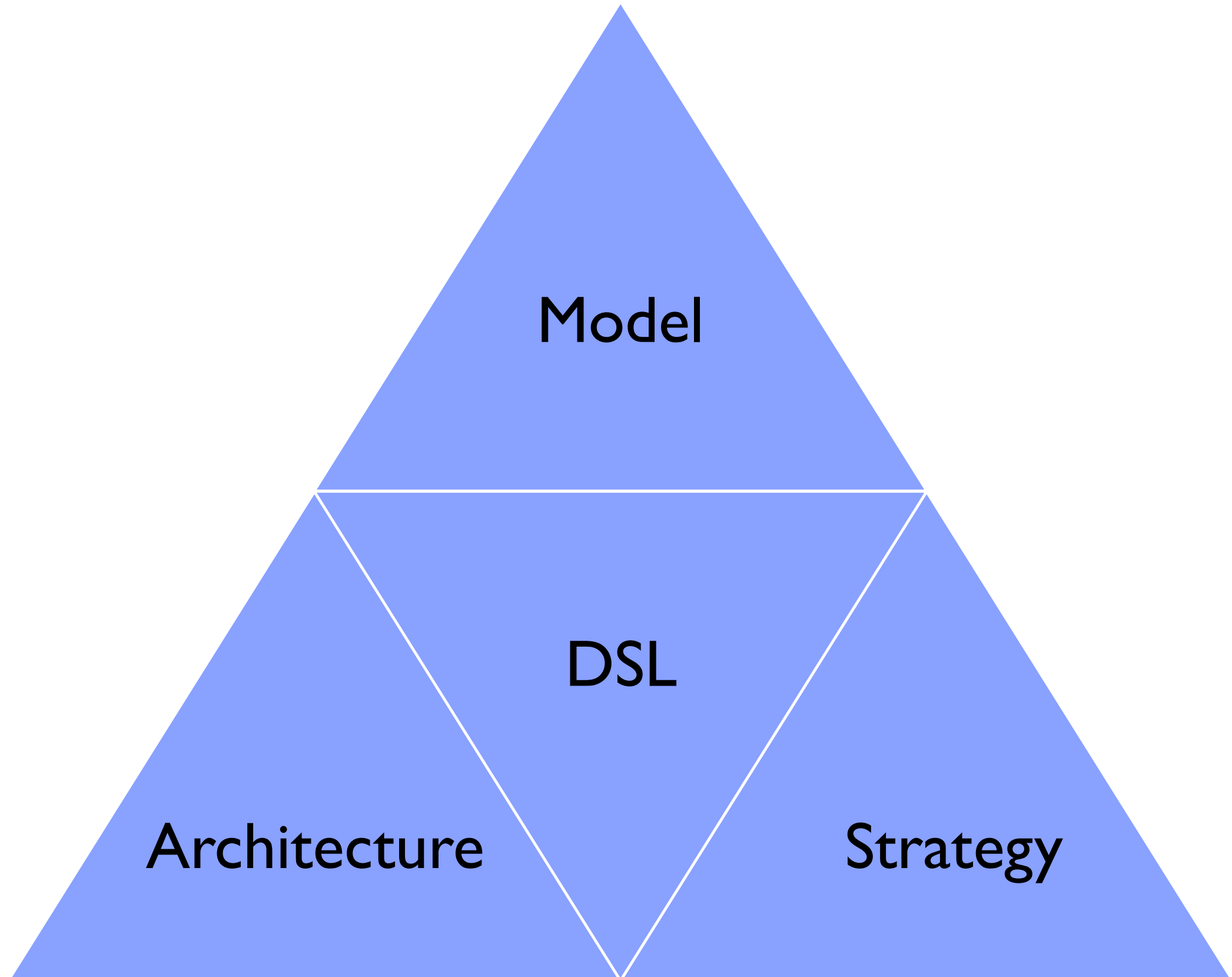


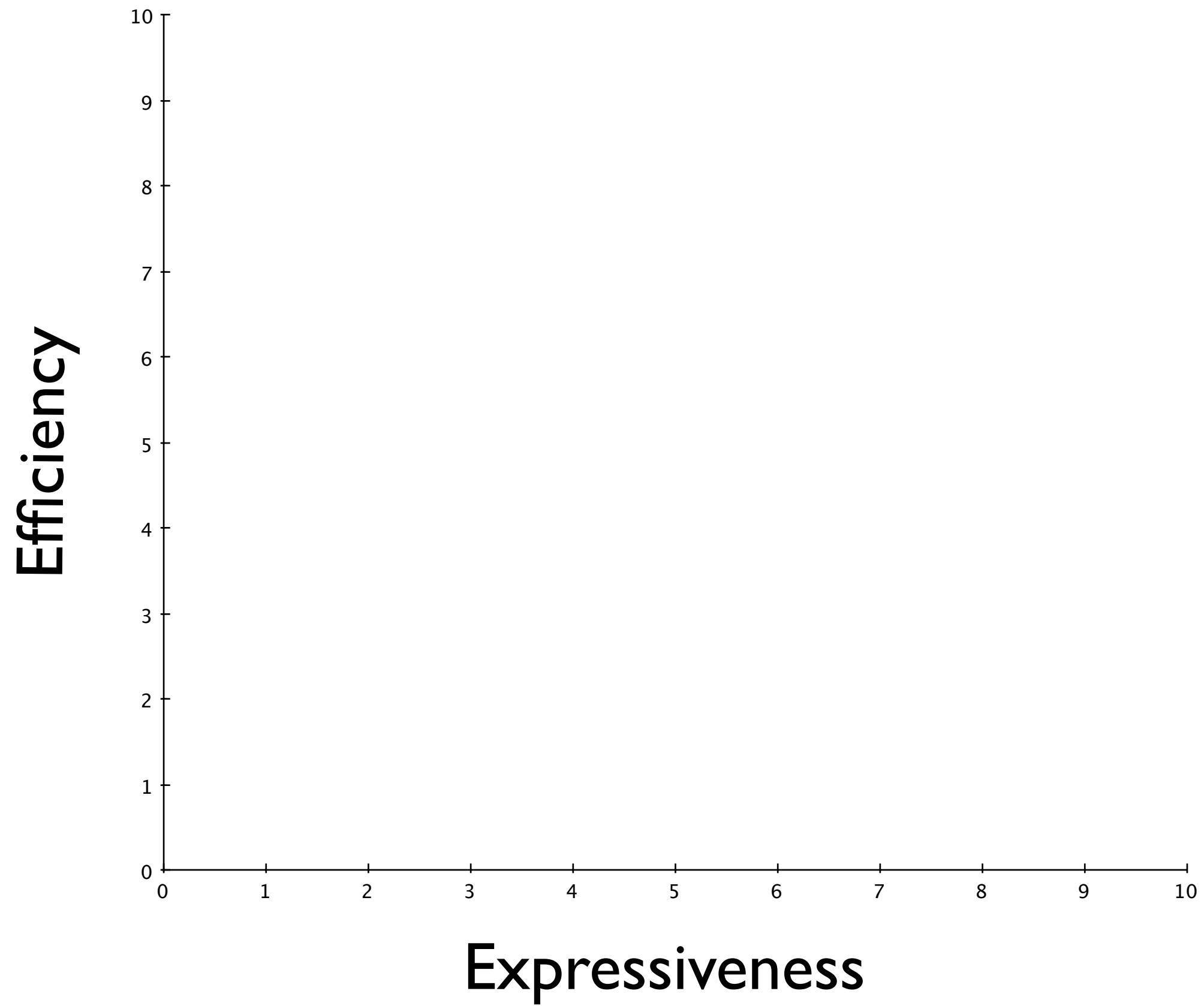


How do we solve both?

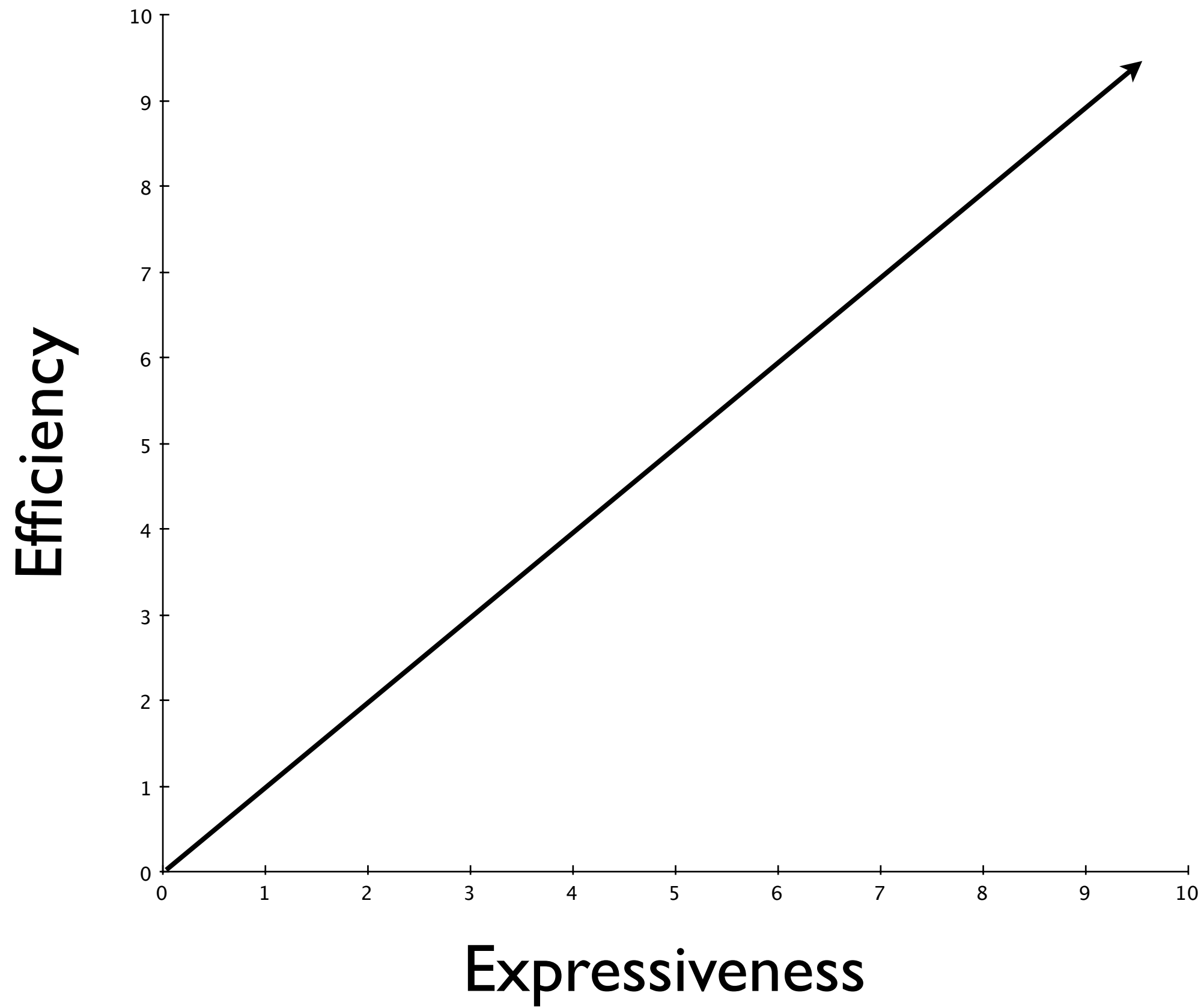




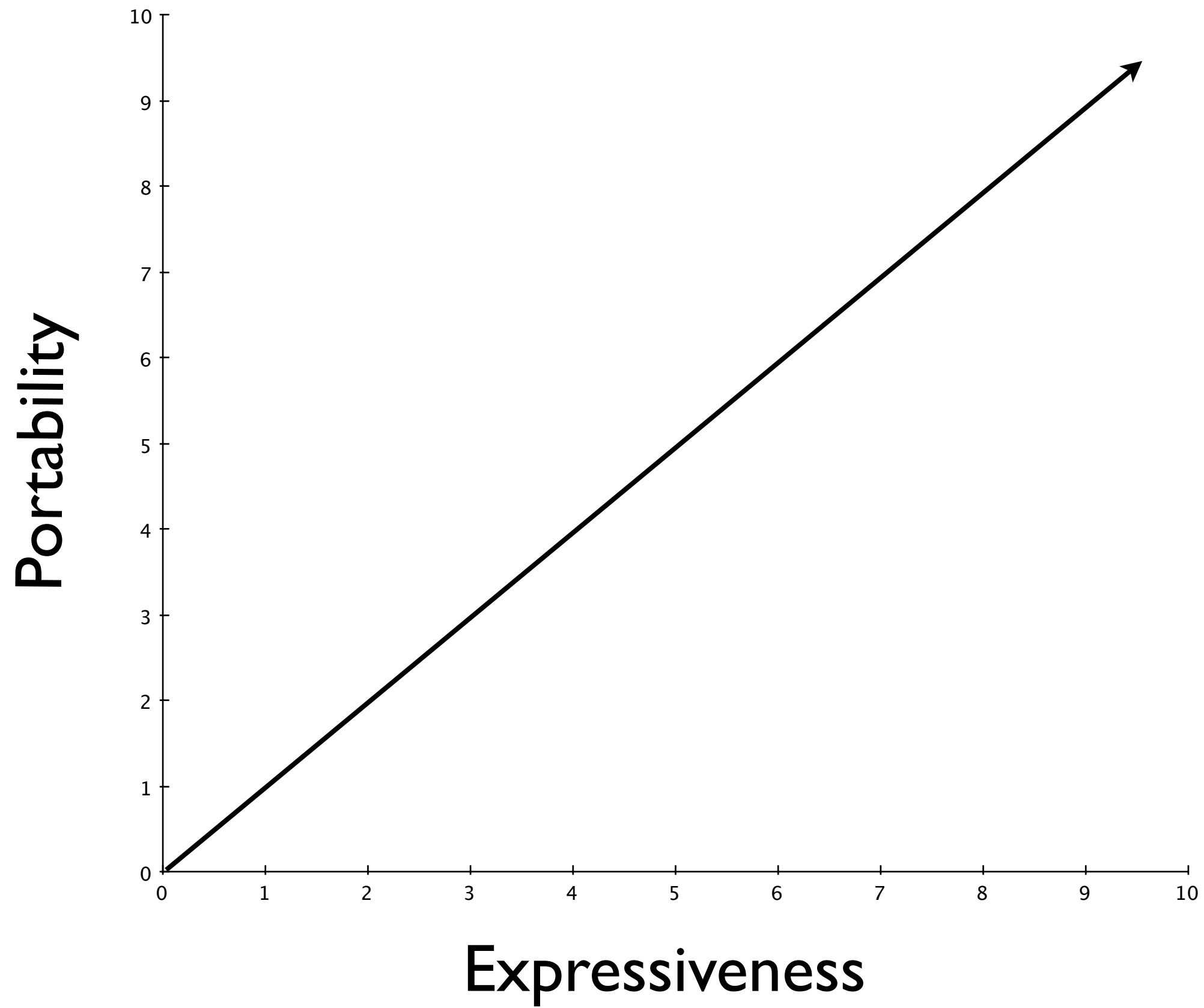












But...

Legacy code matters.

C++

Template

Metaprogramming

Domain Specific Languages

Embedded Domain Specific Languages

Example

$$\vec{c} = \vec{a} + \sin(\vec{b})$$

$c \ll a + \sin(b);$

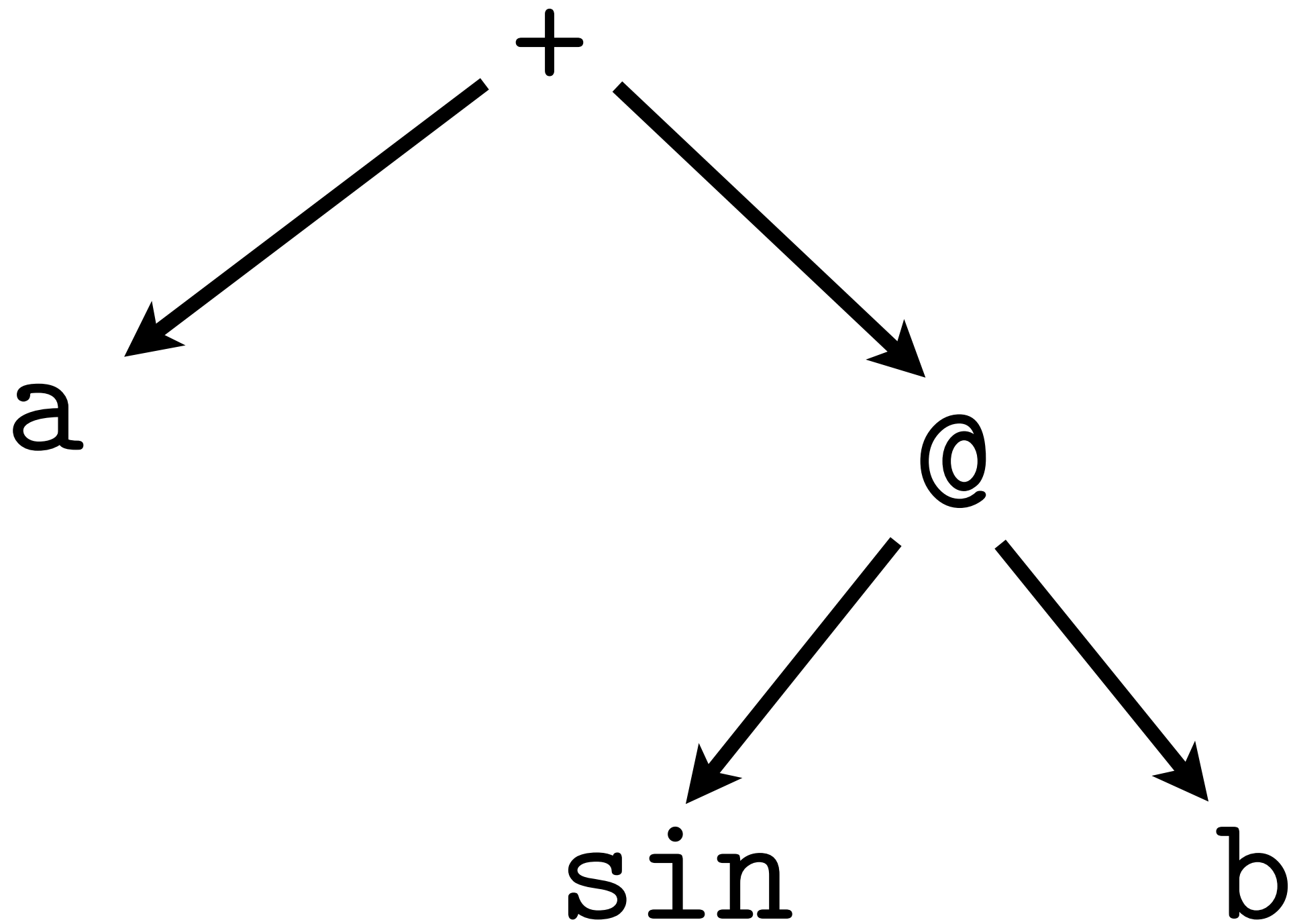
+

a

@

sin

b



$$\vec{c} = \vec{a} + \sin(\vec{b})$$

$c \ll= a + \sin(b); \Rightarrow$

Thread 1 {
Field::const_iterator ia₁ = a₁.begin();
Field::const_iterator ib₁ = b₁.begin();
for(Field::iterator ic₁ = c₁.begin();
 ic₁ != c₁.end();
 ++ic₁, ++ia₁, ++ib₁) {
 *ic₁ = *ia₁ + sin(*ib₁);
 };
:
...
Thread n {
Field::const_iterator ia_n = a_n.begin();
Field::const_iterator ib_n = b_n.begin();
for(Field::iterator ic_n = c_n.begin();
 ic_n != c_n.end();
 ++ic_n, ++ia_n, ++ib_n) {
 *ic_n = *ia_n + sin(*ib_n);
 };

Needed

- Need: Very good with C++
- Want: Template meta-programming
- Want: Comfort with syntax trees

What is lexical analysis?

Characters in;

tokens out.

```
function id(x)
{
    return x ;    // comment
}
```

FUNCTION

```
    id(x)
{
    return x ;    // comment
}
```

FUNCTION

IDENT(id)

(x)

```
{  
  return x ;  // comment  
}
```

FUNCTION

IDENT(id)

LPAR

x)

{

return x ; // comment

}

FUNCTION

IDENT(id)

LPAR

IDENT(x)

)

{

return x ; // comment

}

FUNCTION

IDENT(id)

LPAR

IDENT(x)

RPAR

```
{  
    return x ;    // comment  
}
```


FUNCTION

IDENT(id)

LPAR

IDENT(x)

RPAR

LBRACE

```
return x ; // comment  
}
```

	FUNCTION
	IDENT(id)
	LPAR
	IDENT(x)
	RPAR
x ; // comment	LBRACE
}	RETURN

```
FUNCTION
IDENT(id)
LPAR
IDENT(x)
RPAR
LBACE
; // comment
RETURN
IDENT(x)
```

	FUNCTION
	IDENT(id)
	LPAR
	IDENT(x)
	RPAR
// comment	LBRACE
}	RETURN
	IDENT(x)
	SEMI

FUNCTION

IDENT(id)

LPAR

IDENT(x)

RPAR

LBRACE

RETURN

IDENT(x)

SEMI

}

FUNCTION

IDENT(id)

LPAR

IDENT(x)

RPAR

LBRACE

RETURN

IDENT(x)

SEMI

RBRACE

FUNCTION

LPAR

RPAR

LBRACE

IDENT ()

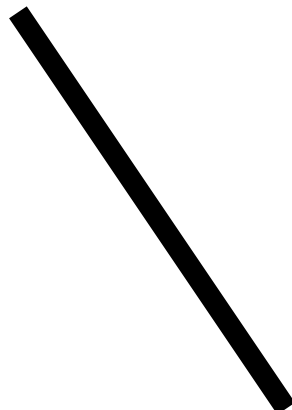
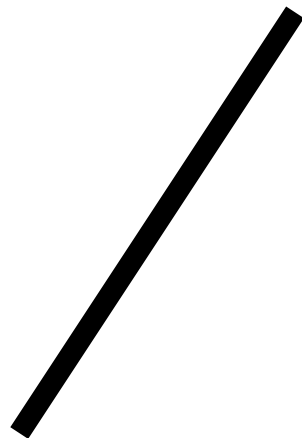
IDENT ()

RETURN

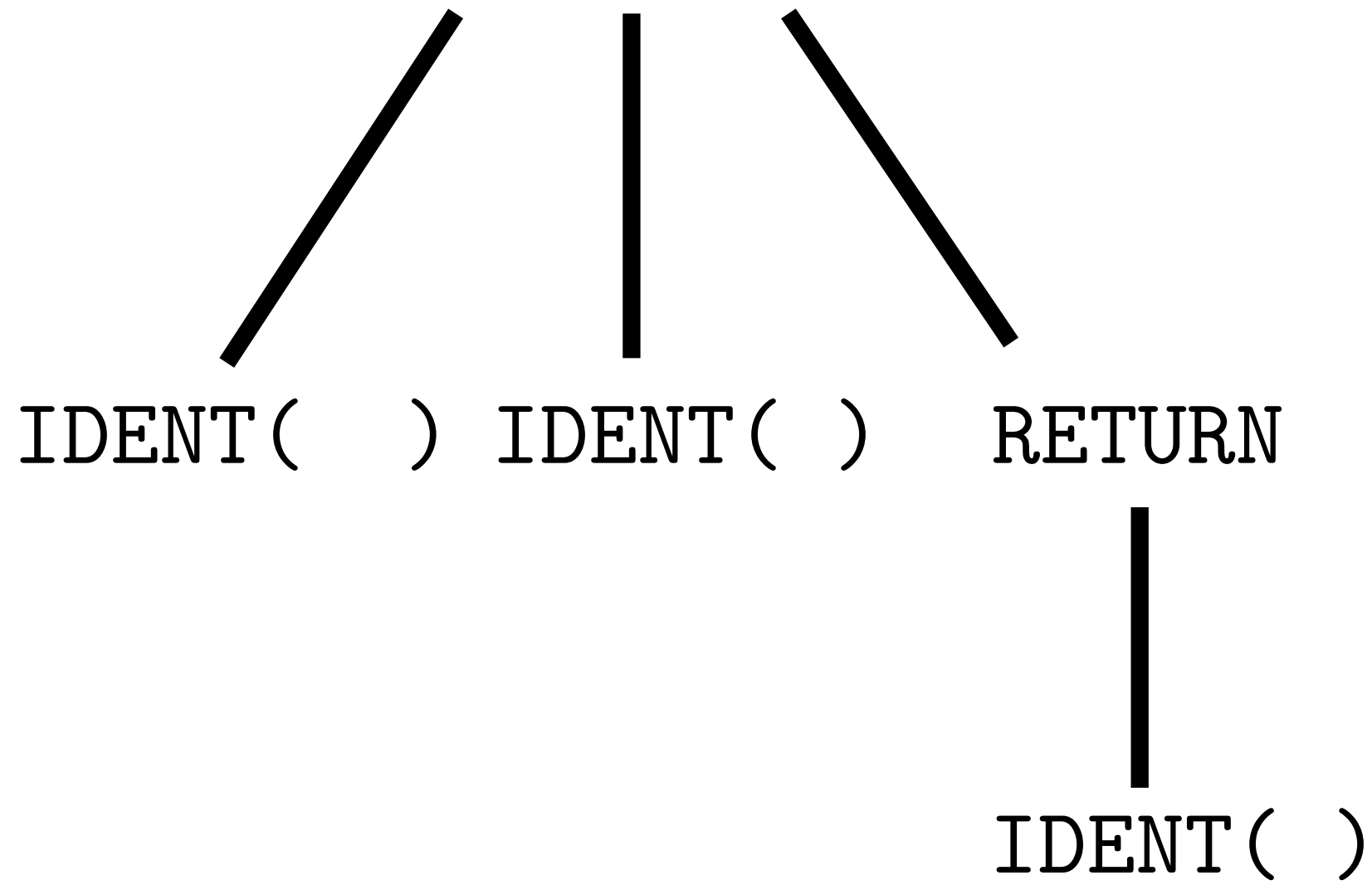
IDENT ()

SEMI

RBRACE



FUNCTION



FUNCTION



How do you define a lexer?

Lexical specification

- Token types
- Whitespace
- Keywords
- Operators
- Comments
- Identifiers
- Punctuation

How do you define tokens?

regex

- Identifiers match $[A-Za-z_][A-Za-z0-9_]*$
- Delimiters match $[(\);]$
- Operators match $[-*+/\wedge=]$
- Integers match $-?([1-9][0-9]*|0)$
- Whitespace is ignored
- Comments match $^\#[^\n]*$

regex is a pattern language.

(a family of them, actually)

regex languages

- math
- unix
- pcre

Regexen match strings.

Minimal regex

Pattern forms

Character: *c*

ex: a matches a

Sequence: $p_1 p_2$

ex: ab matches ab

Choice: $p_1 \mid p_2$

ex: $a \mid b$ matches a

ex: $a \mid b$ matches b

Repetition: p^*

ex: a^* matches a

ex: a^* matches aa

ex: a^* matches aaa

ex: a^* matches

Subpattern: (*p*)

ex: $(ab)^*$ matches ab

ex: $(ab)^*$ matches abab

ex: ab^* matches abbbbbb

Sugared regex

Option: p ?

ex: $a?$ matches a

ex: a? matches

$$p^? = (p \mid)$$

Strict repetition: p^+

ex: a^+ matches a

ex: a^+ !matches

$$p^+ = pp^*$$

Char set: $[c_1 c_2 \dots]$

ex: `[ab]` matches `a`

ex: [ab] matches b

Inverse char set: $[\hat{} c_1 c_2 \dots]$

Ranged char set: $[c_1 - c_2]$

ex: `[a-c]` matches `b`

ex: `[^ab]` matches `c`

Warning: echo \$LANG

Solution: LANG=C

Any character: .

ex: . matches a

ex: . matches c

ex: . matches b

Start of line/string: ^

End of line/string: \$

Repeat: $p\{n\}$

$$\text{ex: } p\{3\} = ppp$$

Bounded repeat: $p\{n, m\}$

$$\text{ex: } p\{2, 4\} = pp \mid ppp \mid pppp$$

Irregular expressions

Submatch: (*p*) and \1

ex: (.)\1 matches aa

Shortest match: p^+ ?

Dialects

Dialect: BRE

- () => \(\)
- { } => \{ \}
- ? => \?
- + => \+
- | => \|

Dialect: BRE, ERE, PCRE

Defined char classes

BRE, ERE	PCRE
<code>[:word:]</code>	<code>\w</code>
<code>[:alpha:]</code>	<code>[A-Za-z]</code>
<code>[:space:]</code>	<code>\s</code>
<code>[:lower:]</code>	<code>[a-z]</code>
<code>[:upper:]</code>	<code>[A-Z]</code>
<code>[:digit:]</code>	<code>\d</code>
<code>[:punct:]</code>	

Define char classes

BRE, ERE	PCRE
<code>\b</code>	<code>\b</code>
<code>[^[:space:]]</code>	<code>\S</code>
<code>[^[:word:]]</code>	<code>\W</code>
<code>[^[:digit:]]</code>	<code>\D</code>

grep: examples and pitfalls

```
grep foo|bar words
```

```
grep 'foo|bar' words
```

```
grep 'foo\|bar' words
```

```
egrep 'foo|bar' words
```

```
grep '.x.n' words
```

```
grep '^ .x.n$' words
```



```
grep '^(.*)\1$' words
```

```
egrep '^(.*)\1$' words
```

```
egrep '98.17.132.45' log
```

```
egrep '\b98\.17\.132\.45\b' log
```

```
egrep '([0-255]\.){3}[0-255]' log
```

\d|1?\d\d|2[0-4]\d|25[0-5]

\d*[02468]

```
perl -e 'while (<STDIN>) {  
    if (/^1?$|^(11+?)\1+$/)  
    { print }  
}'
```


Challenge

- Match dates YYYY MM DD
- Account for leap years
 - If multiple of 400: Yes
 - If multiple of 100 but not 400: No
 - If multiple of 4 but not 100: Yes

Challenge

- Match RFC 3696 email addresses
- Never try this with lives at stake

- Identifiers match $[A-Za-z_][A-Za-z0-9_]*$
- Delimiters match $[(\) ;]$
- Operators match $[-*+/\wedge=]$
- Integers match $-?([1-9][0-9]^*|0)$
- Whitespace is ignored
- Comments match $^\#[\wedge\backslash n]^*$

Python tokens

Integer literals

integer	::=	<u>decimalinteger</u> <u>octinteger</u> <u>hexinteger</u> <u>bininteger</u>
decimalinteger	::=	<u>nonzerodigit</u> <u>digit</u> * "0"+
nonzerodigit	::=	"1" ... "9"
digit	::=	"0" ... "9"
octinteger	::=	"0" ("o" "O") <u>octdigit</u> +
hexinteger	::=	"0" ("x" "X") <u>hexdigit</u> +
bininteger	::=	"0" ("b" "B") <u>bindigit</u> +
octdigit	::=	"0" ... "7"
hexdigit	::=	<u>digit</u> "a" ... "f" "A" ... "F"
bindigit	::=	"0" "1"

Floating point

floatnumber	::=	<u>pointfloat</u> <u>exponentfloat</u>
pointfloat	::=	[<u>intpart</u>] <u>fraction</u> <u>intpart</u> "."
exponentfloat	::=	(<u>intpart</u> <u>pointfloat</u>) <u>exponent</u>
intpart	::=	<u>digit</u> +
fraction	::=	"." <u>digit</u> +
exponent	::=	("e" "E") ["+" "-"] <u>digit</u> +

Operators

+

-

*

**

/

//

%

<<

>>

&

|

^

~

<

>

<=

>=

==

!=

Questions?