# Notes for Polynomial Language Talk

I want a language for polynomials

$$e ::= c \mid x \mid e_0 + e_1 \mid e_0 \cdot e_1$$
$$\mid \text{let } x = e_0 \text{ in } e_1$$

Define top-level functions/programs

$$P ::= \text{def name}(x_0, x_1, \ldots) \; e$$

---

Now, what might I want to do with a program written in this grammar?

- Interpret/Evaluate  ⎤ Semantics
- Compile             ⎦

- "Check" it          ⎤ Safety/
                       ⎦ Types

- take Derivative       ⎤
- take Adjoint Deriv.   ⎥ Transform
                        ⎥ Program
- Partially Evaluate ⎦

- "Simplify"
- Canonicalize
- Optimize

(changing semantics)

Semantics-Preserving Transformations

---

Much of this remains true if we modify the language.

Modifications?
- Add Arrays ?
- Numeric Integration ?
- Probability Distributions
- Dist. Computing Primitives

— ...

---

Interpreting POLY

$$[\![ e \mid \sigma ]\!] = value$$

$e$ — expression
$\sigma$ — environment ($\sigma$ : variable names $\rightarrow$ values)

$[\![ \quad \mid \quad ]\!]$

$$[\![ c \mid \sigma ]\!] = c$$

$$[\![ x \mid \sigma ]\!] = \sigma(x)$$

$$[\![ e_0 + e_1 \mid \sigma ]\!] = [\![ e_0 \mid \sigma ]\!] + [\![ e_1 \mid \sigma ]\!]$$

$$[\![ e_0 \cdot e_1 \mid \sigma ]\!] = [\![ e_0 \mid \sigma ]\!] \cdot [\![ e_1 \mid \sigma ]\!]$$

$$\left[\!\!\left[ \begin{array}{c} \text{let } x = e_0 \\ \text{in } e_1 \end{array} \,\middle|\, \sigma \right]\!\!\right] = \left[\!\!\left[ e_1 \,\middle|\, \sigma\left[ x \mapsto [\![ e_0 \mid \sigma ]\!] \right] \right]\!\!\right]$$

where
$$\sigma\left[ x \mapsto v \right](y) = \begin{cases} v, & \text{if } x = y \\ \sigma(y), & \text{otherwise} \end{cases}$$

---

## Show Code

Compile ...

(let's make a simplifying assumption)

e has the form

$$Le ::= \text{let } x = e \text{ in } Le$$
$$\mid \text{ return } e$$

$$e ::= x \mid c \mid e_0 + e_1 \mid e_0 \cdot e_1$$

That is, all let-bindings are at

the top level.

Then we can compile to a C-function

$$C[\![\ \text{def name}(x_0, x_1, \ldots)\ e\ ]\!]$$

$$= \text{double name}(\\
\quad \text{double } x_0,\\
\quad \text{double } x_1,\\
\quad \ldots\\
)\ \{$$

$$\qquad C[\![\ e\ ]\!]$$

$$\}$$

$$C\left[\!\!\begin{array}{c} \text{let } x = e \\ \text{in } Le \end{array}\!\!\right] = \begin{array}{l} \text{double } x = C[\![\ e\ ]\!]\, ; \\ C[\![\ Le\ ]\!] \end{array}$$

$$C[\![\ \text{return } e\ ]\!] = \text{return } C[\![\ e\ ]\!]\, ;$$

(and you can do the rest)
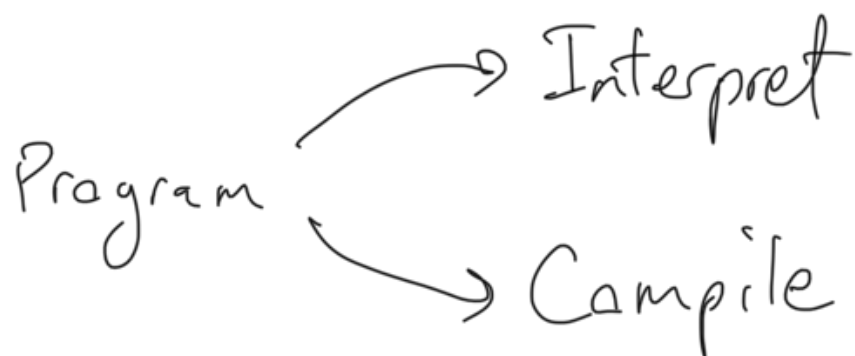
---

# Checks

① What happens if there's a
free variable that isn't bound
during evaluation   or isn't

in the signature during
compilation?

② Program → Interpret
        ↘ Compile

If a program P interprets w/o error
will it compile w/o error?

If a program P compiles w/o error
will it interpret w/o error?

Ideally these are congruent
for all "well-formed" programs

grammatical

well-formed

Enforce via "checking" safety
(here variable binding)
→ e.g. type-checking
   → "effect-checking"
   → parallel safety        etc...
Defined separately from semantics

# Derivatives

"Total" or "forward" derivative

$$f : \underbrace{\mathbb{R} \times \mathbb{R} \times \cdots \times \mathbb{R}}_{n} \to \mathbb{R}$$

$$Df : \underbrace{\mathbb{R} \times \mathbb{R} \times \cdots \times \mathbb{R}}_{n} \times \underbrace{\mathbb{R} \times \mathbb{R} \times \cdots \times \mathbb{R}}_{n} \to \mathbb{R}$$

e.g.
$$f(x, y) = \cdots$$
$$Df(x, y; dx, dy) = \cdots$$

as partials...

$$Df(x, y; dx, dy) = \frac{\partial f}{\partial x}\bigg|_{x,y} \cdot dx + \frac{\partial f}{\partial y}\bigg|_{x,y} \cdot dy$$

from first principles...

$Df(x, y)$ is the closest linear approximation to

$$f(x + dx, y + dy) \approx f(x, y) + Df(x, y; dx, dy)$$

$$D[\![ \text{def name}(x_0, x_1, \cdots) \, e ]\!]$$

$$= \text{def Dname}(x_0, x_1, \cdots, dx_0, dx_1, \cdots)$$

$$D[\![\, e \mid [x_0 \mapsto dx_0, x_1 \mapsto dx_1, \ldots ]\,]\!]$$

$$D[\![\, \begin{array}{l} \text{let } x = e_0 \\ \text{in } e_1 \end{array} \mid \sigma \,]\!] = \begin{array}{l} \text{let } x = e_0 \text{ in} \\ \text{let } dx = D[\![\, e_0 \mid \sigma \,]\!] \\ \text{in } D[\![\, e_1 \mid \sigma[x \mapsto dx] \,]\!] \end{array}$$

This is the Chain Rule

$$D[\![\, c \mid \sigma \,]\!] = \emptyset$$

$$D[\![\, x \mid \sigma \,]\!] = \emptyset \text{ if } x \notin \sigma$$

$$D[\![\, x \mid \sigma \,]\!] = \sigma(x)$$

$$D[\![\, e_0 + e_1 \mid \sigma \,]\!] = D[\![\, e_0 \mid \sigma \,]\!] + D[\![\, e_1 \mid \sigma \,]\!]$$

$$D[\![\, e_0 \cdot e_1 \mid \sigma \,]\!] = D[\![\, e_0 \mid \sigma \,]\!] \cdot e_1 +$$
$$e_0 \cdot D[\![\, e_1 \mid \sigma \,]\!]$$

Leibniz Product Rule

---

"Adjoint Derivative"
   e.g. Gradient
↳ aka. "Reverse Mode"

example: $f(x, \ldots)$

$$f : \mathbb{R} \times \mathbb{R} \longrightarrow \mathbb{R}$$

$$Df : \mathbb{R} \times \mathbb{R} \times \underbrace{\mathbb{R} \times \mathbb{R} \longrightarrow \mathbb{R}}_{\text{linear}}$$

$$D^T f : \mathbb{R} \times \mathbb{R} \times \underbrace{\mathbb{R} \longrightarrow \mathbb{R} \times \mathbb{R}}_{\text{linear}}$$

(same as "matrix transpose")

Algebraic definition by a universal property...

Let $g : A \to B$ linear, then

$g^T : B \to A$ is the unique linear function s.t.

$$\forall x \in A : \forall y \in B :$$

$$\langle y \mid g(x) \rangle = \langle g^T(y) \mid x \rangle$$

$(\langle \cdot \mid \cdot \rangle$ is "inner product"
aka. dot-product$)$

Note:

for $f(x,y)$ $f^T$ returns

a _pair_ but we don't have pairs in our language.

Therefore our language is not _closed_ under adjoint derivatives.

We will ignore this, but note that there is no magic reason that when you investigate a language, everything will work out. That's why this is research.

## Partial Evaluation

We defined the "total derivative" as the closest linear approximation, but what if we want a plain old rise-over-run/slope derivative?

Let $f: \mathbb{R} \to \mathbb{R}$

$f(x)$.

Then $\dfrac{df}{dx} = Df(x; 1)$

$$b/c \quad Df(x, dx) = \frac{df}{dx}\Big|_x \cdot dx$$

For a function definition

$$\text{def name}(x_0, x_1, \ldots) \; e \; ] \quad i$$

the partial evaluation w.r.t. $x_k = V$

is

$$\text{def PEname}(x_0, \ldots, x_{k-1}, x_{k+1}, \ldots) \Big]$$
$$\text{let } x_k = V \text{ in } e$$

<u>Note</u>: this definition of partial evaluation is functionally correct, but does not optimize the implementation how we would normally expect.

---

The preceding program transformations change the signature of a function, so we know they aren't "semantics preserving." By contrast, the remaining transformations will be. Therefore it is useful to define a notion of program equality. The

kind of equality we care about (between expressions) is usually called **Observational Equality**.

$$e_0 \cong e_1 \text{ iff.}$$

$$\forall \sigma : [\![ e_0 \mid \sigma ]\!] = [\![ e_1 \mid \sigma ]\!]$$

where $\sigma$ is such that both sides are well defined.
(i.e. $\sigma$ binds all free variables)

> Observational Eq. is well-defined for general $\lambda$-calculi but we can simplify in our context.

---

Some Random "Simplification" Rules

$$0 + e \cong e$$

$$0 \cdot e \cong 0$$

$$1 \cdot e \cong e$$

$$e + e \cong 2 \cdot e$$

$$\ldots$$

(Note: these rules can be justified by appealing to the semantics)

When will I have "enough" simplification rules? Is that possible?

# Canonical Forms

Let $E$ be the set of well-formed expressions in POLY. Then a canonicalization procedure/function $\gamma : E \to E$ is a function s.t. $\forall e : e \cong \gamma(e)$

(syntactic equality) $\downarrow$

and $\forall e_1, e_2 : (e_1 \cong e_2) \iff \gamma(e_1) = \gamma(e_2)$

That is all equivalent expressions canonicalize into syntactically identical canonical forms.

## Observe:

If we enumerate every rewrite rule used in a canonicalization procedure, then those rewrites effectively __axiomatize__ observational equality.

__Claim__ : This will be possible for POLY, but not for general Turing-complete languages.

The existence of a computable canonicalization function for a language trivially yields a decision procedure for observational equality of terms in that language. However, this is an undecidable problem for Turing-Complete languages

(by Rice's Theorem; see any Automata Theory textbook)

## Canonicalization of POLY

Any polynomial can be uniquely expressed on the "monomial basis."

For instance

$$(x-y) \cdot (x \cdot z + 2) = x^2 z - xyz + 2x - 2y$$

Here is a general procedure.

(1) Substitute all let-bindings

let $x = e_0$ in $e_1 \cong [x \mapsto e_0] e_1$

(2) Distribute all mult. over all sums

$$e_0 \cdot (e_1 + e_2) \overset{\leadsto}{=\!=} e_0 \cdot e_1 + e_0 \cdot e_2$$

sim. for $(e_0 + e_1) \cdot e_2$

(3) Canonicalize monomials

$$e_0 \cdot (e_1 \cdot e_2) \overset{\leadsto}{=\!=} (e_0 \cdot e_1) \cdot e_2$$

$$C_0 \cdot C_1 \overset{\leadsto}{=\!=} c' \leftarrow \text{the computed product}$$

$$e_0 \cdot e_1 \overset{\leadsto}{=\!=} e_1 \cdot e_0$$

↳ apply in conjunction with associativity until all constants are leading and all variables occur in lexicographic order. If there is no leading constant, make 1 the leading constant.

(4) Canonicalize order of summation

$$e_0 + (e_1 + e_2) \overset{\leadsto}{=\!=} (e_0 + e_1) + e_2$$

$$e_0 + e_1 \overset{\leadsto}{=\!=} e_1 + e_2$$

↳ Commute and associate in lexicographic order, with all constants leading

(5) Collapse summation

$$C_0 + C_1 \overset{\leadsto}{=\!=} c' \leftarrow (\text{computed})$$

$$c_0 \cdot e + c_1 \cdot e \stackrel{\sim}{=} c' \cdot e$$

(6) Constant elimination

$$0 \cdot e \stackrel{\sim}{=} 0$$
$$1 \cdot e \stackrel{\sim}{=} e$$
$$0 + e \stackrel{\sim}{=} e$$

This procedure $\gamma$ satisfies

$$\gamma(e) \stackrel{\sim}{=} e \quad \text{by showing that}$$
each constituent rewrite preserves equality.

However, we need some special other argument to show that if

$$\gamma(e_1) \neq \gamma(e_2), \text{ then}$$

$$e_1 \neq e_2 .$$

For polynomials this is simply the fundamental theorem of algebra.

___

We now have a complete set of rewrite rules. But if we want to "simplify" or "optimize"

which do we apply? ¿

In general, substitution and/or distribution can lead to exponentially large terms.

Note: We are now implicitly assuming a <u>Cost-model</u> based on expression size. In more complex languages we may want other cost-models — perhaps even many for the same language.

Rules w/ cast

$$\text{let } x = e_0 \text{ in } e_1 \;\overset{\sim}{\equiv}\; [x \mapsto e_0] e_1$$

$$e_0 \cdot (e_1 + e_2) \;\overset{\sim}{\equiv}\; e_0 \cdot e_1 + e_0 \cdot e_2$$

$$e_0 \cdot (e_1 \cdot e_2) \;\overset{\sim}{\equiv}\; (e_0 \cdot e_1) \cdot e_2$$

$$e_0 \cdot e_1 \;\overset{\sim}{\equiv}\; e_1 \cdot e_0$$

$$e_0 + (e_1 + e_2) \;\overset{\sim}{\equiv}\; (e_0 + e_1) + e_2$$

$$e_0 + e_1 \;\overset{\sim}{\equiv}\; e_1 + e_0$$

$$c_0 \cdot c_1 \;\overset{\sim}{\equiv}\; c' \quad \leftarrow \text{computed}$$

$$c_0 + c_1 \;\overset{\sim}{\equiv}\; c' \quad \leftarrow \text{computed}$$

(leading constants are handleable by distribution + computation)

$$0 \cdot e \cong 0$$
$$1 \cdot e \cong e$$
$$0 + e \cong e$$

Given this analysis, a simplification or optimization procedure would run only the cost-reducing rules, possibly augmented by the normalization orderings of the cost-neutral rules.

However, we can do better by running cost <u>negative</u> rules in <u>reverse</u>.

Running let-substitution in reverse is just <u>common subexpression elimination</u>.

What about distribution?

Consider a polynomial in one variable

$$a_0 + a_1 x + a_2 x^2 + a_3 x^3$$

Then a _Horner Scheme_ is the factorization

$$a_0 + x \cdot (a_1 + x \cdot (a_2 + x \cdot a_3))$$

A multi-variate Horner Scheme is a (usually greedy) factorization of "common terms" from a sum of products.

This is an NP-Hard search problem.

## Note: Optimization has largely reduced to a known problem ...

## Recap

We defined our POLY language

$$e ::= c \mid x \mid e_0 + e_1 \mid e_0 \cdot e_1$$
$$\mid \text{let } x = e_0 \text{ in } e_1$$

To handle adjoints, we need pairs

$$e ::= \cdots \mid (e_0, e_1) \mid \pi_0 e \mid \pi_1 e$$

ATL adds

$$e ::= \cdots \mid \boxplus_{i=0}^{n} e \mid \sum_{i=0}^{n} e \mid e[i]$$

$$\mid [p] \cdot e$$

$$p ::= i = j \mid p_0 \wedge p_1$$

$$\mid \text{\textasciitilde\textasciitilde} \text{ affine expressions } \text{\textasciitilde\textasciitilde}$$
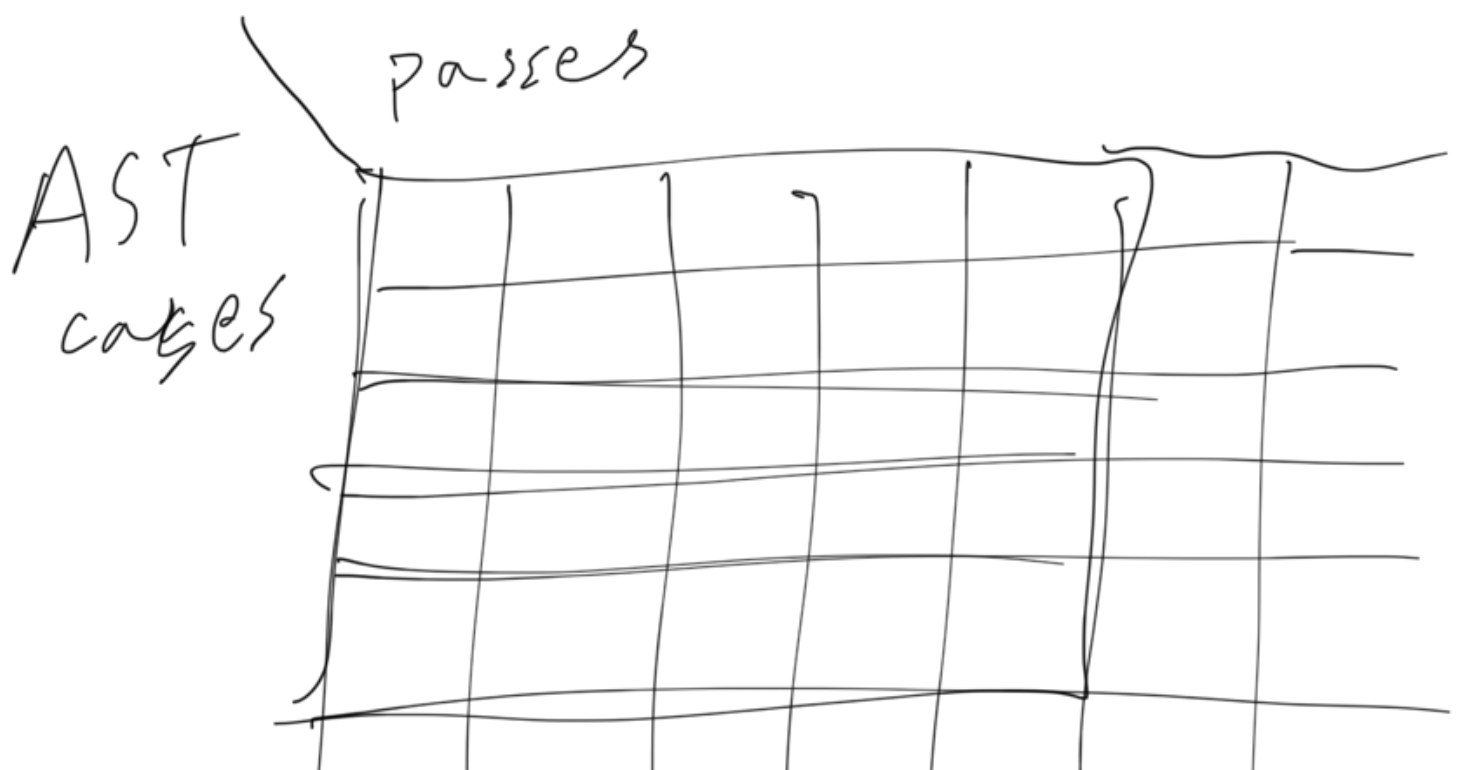
$$\mid R(i_0, i_1, \ldots)$$

What about integrals?

---

# The <u>expression problem</u>

is the combinatorial code growth
induced by growing compilers

passes

AST
cases

↱ Recap