# Model Simplification



Prof. Fabrício Olivetti de França

Federal University of ABC

15 June, 2024



# **Simplifying Algebraic Expressions**



Rewriting an algebraic expression by applying equivalence rules to generate an equivalent expression that is more efficient and compact.



$$x + y \equiv y + x$$
$$2x \equiv x + x$$
$$\tanh x \equiv \frac{e^{2x} - 1}{e^{2x} + 1}$$

There is no standard measure of efficiency and compactness for an algebraic expression as this depends on the application:

- Smaller expressions
- Computationally more efficient
- Less repeated use of the same variable
- Less numerical constants
- etc.

For example:

- x + x + x could be much faster to calculate than 3x in older architectures.
- 3x is more precise if we are using interval arithmetic.
- (x + 1.5y)/(1 + 3x) can be more interpretable than (2x + 3y)/(2 + 6x) (if these parameters can still fit the data).

This simplification process depends on a set of equivalence rules that describes whether two expressions are equivalent. Let's see a couple more examples:

$$x + y \equiv y + x$$
$$xy \equiv yx$$
$$2x \equiv x + x$$
$$x + (y + z) \equiv (x + y) + z$$
$$x(y + z) \equiv xy + xz$$
$$(xy)/z \equiv x(y/z)$$

Finding the optimal expression is not as simple as applying these rules sequentially. The order of the application matters, enumerating all possible orders leads to a combinatorial explosion.

If we have N rules and we can only apply them once, we have N! ways of applying these rules.

Of course, not every rule is appliable at every moment and we may need to apply a rule more than once.

## This is a well known problem in compiler optimization called **Phase Or**dering Problem:



Given a measure of goodness for each rule at a given expression, if two rules give the same local benefit, then both can lead to sub-optimal results.

Let's see this problem using algebraic expressions (2x)/2 and the following rules:

$$2\alpha \equiv \alpha << 1$$
$$\alpha\beta \equiv \beta\alpha$$
$$1\alpha \equiv \alpha$$
$$\alpha/\alpha \equiv 1$$
$$(\alpha\beta)/\gamma \equiv \alpha(\beta/\gamma)$$

$$2\alpha \equiv \alpha << 1$$
$$\alpha\beta \equiv \beta\alpha$$
$$1\alpha \equiv \alpha$$
$$\alpha/\alpha \equiv 1$$
$$(\alpha\beta)/\gamma \equiv \alpha(\beta/\gamma)$$

 $\alpha = x$ (2x)/2 = (x << 1)/2

$$2\alpha \equiv \alpha << 1$$
$$\alpha\beta \equiv \beta\alpha$$
$$1\alpha \equiv \alpha$$
$$\alpha/\alpha \equiv 1$$
$$(\alpha\beta)/\gamma \equiv \alpha(\beta/\gamma)$$

$$2\alpha \equiv \alpha << 1$$
$$\alpha\beta \equiv \beta\alpha$$
$$\boxed{1\alpha \equiv \alpha}$$
$$\alpha/\alpha \equiv 1$$
$$(\alpha\beta)/\gamma \equiv \alpha(\beta/\gamma)$$

$$2\alpha \equiv \alpha << 1$$
$$\alpha\beta \equiv \beta\alpha$$
$$1\alpha \equiv \alpha$$
$$\alpha/\alpha \equiv 1$$
$$(\alpha\beta)/\gamma \equiv \alpha(\beta/\gamma)$$

$$2\alpha \equiv \alpha << 1$$
$$\alpha\beta \equiv \beta\alpha$$
$$1\alpha \equiv \alpha$$
$$\alpha/\alpha \equiv 1$$
$$(\alpha\beta)/\gamma \equiv \alpha (\beta/\gamma)$$

That was quick and sub-optimal! Let's try again:

$$\alpha\beta \equiv \beta\alpha$$
$$(\alpha\beta)/\gamma \equiv \alpha(\beta/\gamma)$$
$$\alpha/\alpha \equiv 1$$
$$1\alpha \equiv \alpha$$
$$2\alpha \equiv \alpha << 1$$

$$\alpha = 2$$
$$\beta = x$$
$$(2x)/2 = (x2)/2$$

$$\alpha\beta \equiv \beta\alpha$$

$$(\alpha\beta)/\gamma \equiv \alpha(\beta/\gamma)$$

$$\alpha/\alpha \equiv 1$$

$$1\alpha \equiv \alpha$$

$$2\alpha \equiv \alpha << 1$$

$$\alpha = x$$
  

$$\beta = 2$$
  

$$\gamma = 2$$
  

$$(2x)/2 = (x2)/2$$
  

$$(x2)/2 = x(2/2)$$

$$\alpha\beta \equiv \beta\alpha$$
$$(\alpha\beta)/\gamma \equiv \alpha(\beta/\gamma)$$
$$\alpha/\alpha \equiv 1$$
$$1\alpha \equiv \alpha$$
$$2\alpha \equiv \alpha << 1$$

$$\alpha = 2$$

$$(2x)/2 = (x2)/2$$

$$(x2)/2 = x(2/2)$$

$$x(2/2) = x1$$

$$\alpha\beta \equiv \beta\alpha$$
$$(\alpha\beta)/\gamma \equiv \alpha(\beta/\gamma)$$
$$\alpha/\alpha \equiv 1$$
$$1\alpha \equiv \alpha$$
$$2\alpha \equiv \alpha << 1$$

$$(2x)/2 = (x2)/2$$
  
 $(x2)/2 = x(2/2)$   
 $x(2/2) = x1$ 

So close! But if we go through the rules again we will reach the optimal result!

$$\alpha\beta \equiv \beta\alpha$$
$$(\alpha\beta)/\gamma \equiv \alpha(\beta/\gamma)$$
$$\alpha/\alpha \equiv 1$$
$$1\alpha \equiv \alpha$$
$$2\alpha \equiv \alpha << 1$$

$$(2x)/2 = (x2)/2$$
  
 $(x2)/2 = x(2/2)$   
 $x(2/2) = x1$ 

$$\alpha\beta \equiv \beta\alpha$$
$$(\alpha\beta)/\gamma \equiv \alpha(\beta/\gamma)$$
$$\alpha/\alpha \equiv 1$$
$$1\alpha \equiv \alpha$$
$$2\alpha \equiv \alpha << 1$$

$$\alpha = x$$
$$\beta = 1$$

$$(2x)/2 = (x2)/2$$
  
 $(x2)/2 = x(2/2)$   
 $x(2/2) = x1$   
 $x1 = 1x$ 

#### Good!

$$\alpha\beta \equiv \beta\alpha$$
$$(\alpha\beta)/\gamma \equiv \alpha(\beta/\gamma)$$
$$\alpha/\alpha \equiv 1$$
$$1\alpha \equiv \alpha$$
$$2\alpha \equiv \alpha << 1$$

$$\alpha = x$$

$$(2x)/2 = (x2)/2$$

$$(x2)/2 = x(2/2)$$

$$x(2/2) = x1$$

$$x1 = 1x$$

$$1x = x$$

One reason for the phase ordering problem is that the program transformation is destructive.

Whenever we apply a rule, we *forget* about the past equivalent expressions. So, whenever we take the wrong *path*, we cannot go back.

## **Equality Saturation**

In <sup>1</sup> Tate introduced the idea of **Equality Saturation** proposing a method for non-destructive rewriting of the original expression.

The main idea is that all the rules are applied in parallel and the resulting expressions are kept in a compact data structure now known as **e-graph**.

<sup>1</sup>Tate, R., et al, "Equality Saturation: a New Approach to Optimization," in Logical Methods in Computer Science, 2011.

This technique became popular after a fast and customizable implementation called  $egg^2$ . Egg is implemented in Rust and enabled the user to experiment with equality saturation in different applications with the requirements of:

- Defining a Language
- Writing the rules

<sup>&</sup>lt;sup>2</sup>Willsey, M., et al. "egg: Fast and extensible equality saturation," in Proceedings of the ACM on Programming Languages, vol. 5, no. POPL, pp. 1–29, 2021

An e-graph<sup>3</sup> is composed of:

- a set of e-classes (dashed) with a non-empty sets of e-nodes (solid line).
- a set of e-nodes representing a symbol of our language
- a set of edges connecting e-nodes to e-classes.

<sup>&</sup>lt;sup>3</sup>Nelson, Charles Gregory. Techniques for program verification. Stanford University, 1980.



**Figure 1:** (a) Original expression; (b) after applying the rule  $1 * x \rightarrow x$ ; (c) after applying the rule  $x + x \rightarrow 2 * x$ .

Notice how the children of an e-node points to the e-classes that contain one or more equivalent sub-expression. For example, we can go from + to 1x or x or  $1 \times 1 \times x$ .



At a certain point, whenever we apply the rules the graph will not change! This means we reached the fixed point and the graph is saturated.

The saturated graph represents every equivalent expression reachable by the specific set of rules.



# egg: Fast and Extensible Equality Saturation

```
def equality_saturation(expr, rewrites):
        egraph = build_egraph(expr)
2
        while !egraph.is_saturated_or_timeout():
3
            matches = []
4
5
            # read-only
6
            for rw in rewrites:
7
                for (subst, eclass) in egraph.search(rw.lhs):
8
                    matches.append((rw, subst, eclass))
9
            # write-only
            for (rw, subst, eclass) in matches:
12
                new_eclass = egraph.add(rw.rhs.subst(subst))
13
                egraph.union(eclass, new_eclass)
14
            # restore invariants: merge congruent e-classes
16
            egraph.rebuild()
17
18
        return egraph.extract_best()
19
```

In the following, let's use as an illustrative example the expression (2/x)(x + x) and the rules:

 $\alpha + \alpha \to 2\alpha$  $(\alpha/\beta)\gamma \to \alpha(\gamma/\beta)$  $(\alpha\beta)/\gamma \to \alpha(\beta/\gamma)$  $\alpha/\alpha \to 1$  $\alpha \cdot 1 \to \alpha$ 

The representation of an e-node follows the fixed point of the representation of the program or expression.

In mathematics, a **fixed point** x of a function f, also known as **invariant point** is a value x such that f(x) = x. The **least fixed point** is the smallest fixed point x among all fixed points of a function.

1

As a data structure we can represent a fixed point as:

data Fix f = Fix (f (Fix f))

which allows us to implement generic algorithms for folding and unfolding data structures.
A natural way to represent an expression tree is using a recursive type:

data Expr	= Const Double
	Var String
	Add Expr Expr
	Sub Expr Expr
	Mul Expr Expr
	Div Expr Expr

2

7 8 We can represent the same structure as its least fixed point with:

data	Expr	а	=	Cons	st	Double		
	-			Var	St	ring		
				Add	а	a		
				Sub	а	a		
				Mul	а	a		
				Div	а	a		
type	type FixExpr = Fix Expr							

Notice that this structure allows a recursive pattern by replacing the type parameter a with FixExpr.

This structure also give us many conveniences, we can represent and expression x + 2 as:

<sup>1</sup> Fix (Add (Fix (Var "x") (Fix (Const 2))) :: FixExpr

which types check. But also, we can have represent the e-node Add point to the e-classes 3 and 1:

1 Add 3 1 :: Expr Int

or a representation of the pattern \_ + \_ where each \_ matches anything:

Add () () :: Expr ()

In dynamically typed language this can be expressed as structures with branches that can assume different types as it seems fit.

```
class Add:
        left: Any
2
        right: Any
3
4
    class Mul:
5
        left: Any
6
        right: Any
7
8
    class Var:
9
        name: String
10
11
    class Const:
12
        val: float
13
14
    Expr = Add | Mul | Var | Const
15
```

This definition allows us to represent our expression as an expression tree where each children of a node is another tree or an integer representing the e-class id it points to or a bottom symbol.

We will see next how this flexibility is useful for the implementation of egg.

Starting with our exression in a tree-like representation:



Visiting the nodes in a post-order traversal...



We create a new e-node assiging it a new e-class id. This information is stored in a hashcon, where the key is the node and the value is the corresponding e-class id.



### We create an e-node storing it into a new e-class and assign an e-class id.



Any internal node will have their children replace by the corresponding eclass ids. We keep a hash table of this node representation to the corresponding assigned id.



Whenever we reach a node contained in the hash table, we just retrieve the e-class id it belongs to.



This will ensure that the graph is compact and the equivalent nodes are grouped together.



#### With this representation, it becomes easier to search the patterns for a match.



The corresponding e-graph of the expression (2/x)(x+x) is



### By applying the rule $x + x \rightarrow 2x$ we get





By applying the rule  $(x/y)z \rightarrow x(z/y)$  we get

After this iteration, we must merge the e-nodes that are equivalent to the same e-class:



Notice how following any path from a given e-class will return equivalent expressions.

After applying every rule, we start again ignoring the already performed rewritings. Now we can match  $(xy)/z \rightarrow x(y/z)$ :



## And now $x/x \to 1$ :





The pattern matching procedure starts with the creation of a database of patterns. This database is composed of a hash table where the keys are each node with their children replaced by a unit value, and the values are a sequence of Tries linking to the e-class of id of each element.



For example, the entry for \* contains two e-classes: 1, 3. The first child of e-class 1 can be either 2 or 5 depending of what e-node is chosen. This follows up until for every child of those e-nodes.



If we want to match the rule ("x" / "y") \* "z", we first retrieve the patterns \_ \* \_ from the hash table, finding that it is found at e-classes 1 and 3.



Next, we retrieve the left child of these e-classes, generating the set {2, 5}.



Finally, we retrieve the entry \_ / \_ from the hash table, finding the e-classes {2,4}. Intersecting both sets  $\{2,5\} \cap \{2,4\} = \{2\}$  we find the single match.



After reaching saturation, we can extract the optimal expressions by following a greedy heuristic or any graph traversal algorithm we find suitable.



For example, if we say that + has cost 2 and \*, / have cost 3, and terminals have cost 1. By greedly following the post-order traversal, we would assign a partial result for each e-class as the minimum result from all e-nodes.



#### In our example, we would find 2 \* 2 as the best expression.



Other interesting features from *egg* is the constant folding, which allows us to store facts about each e-class. In our example, we can store the constant that this particular e-class can be evaluated to, if any.



In this situation, we could simply add the e-node 4 into the root e-class, finding an even smaller expression.



 "General characteristic of equality saturation: either a successful rewrite sequence is found relatively quickly, or, computational costs explode" <sup>4</sup>

<sup>&</sup>lt;sup>4</sup>Thomas Kœhler, Andrés Goens, Siddharth Bhat, Tobias Grosser, Phil Trinder, and Michel Steuwer. 2024. Guided Equality Saturation. Proc. ACM Program. Lang. 8, POPL, Article 58 (January 2024), 32 pages. https://doi.org/10.1145/3632900

Take as an example the rules  $(\alpha + \beta) + \gamma \Rightarrow \alpha + (\beta + \gamma), \alpha + 0 \Rightarrow \alpha$ ,  $(-\alpha) + \alpha \Rightarrow 0$ . If we have the expression (x + (-y)) + y, it will generate:

$$(x + (-y)) + y \Rightarrow x + ((-y) + y)$$
  
$$\Rightarrow (x + (-y)) + y) + (-y + y)$$
  
$$\Rightarrow \dots$$

# Nonterminating rules



- egraphs-good website
- egg
- hegg (Haskell)
- Metatheory.jl (Julia)
- egg presentation

# **Application in Symbolic Regression**

Symbolic Regression algorithms are prone to overparametrization<sup>5</sup>:

$$f(\mathbf{x}, \theta) = \theta_1 \exp\left(\theta_2 x_1 + \theta_3\right)$$

- The parameters can assume different values for the same data
- Numerical issues and slow convergence of optimization
- Larger search space for memetic approaches
- Interpretation is hindered

<sup>&</sup>lt;sup>5</sup>de Franca, Fabricio Olivetti, and Gabriel Kronberger. "Reducing Overparameterization of Symbolic Regression Models with Equality Saturation." Proceedings of the Genetic and Evolutionary Computation Conference. 2023.
- Is overparametrization really a problem?
- Can EqSat help to reduce overparametrization in SR?
  - How well does it compare with simpler alternatives? (e.g., Sympy)
- Does EqSat always reduce to the optimal number of parameters?
- How fast is EqSat?

## Is this a real problem?

- · Run MOO version of Operon once for every Feynman dataset
- Stored the Pareto front with a total of  $183\,491$  models



## How much can EqSat help?

• 30 independent runs of a set of SR algorithms to:

$$f_1(x,y) = \frac{1}{1+x^{-4}} + \frac{1}{1+y^{-4}}$$
(1)  
$$f_2(x,y) = \frac{e^{-(x-1)^2}}{1.2 + (y-2.5)^2}$$
(2)

- Simplify with:
  - EqSat
  - Sympy
  - Sympy + EqSat

- SR algorithms:
  - Bingo
  - EPLEX
  - GP-GOMEA
  - Operon
  - PySR
  - SBP

## How much can EqSat help?



## % simplified expressions with # parameters equal to rank (left) or with at most one extra (right).

Algorithm	Pagie-1	Kotanchek	Algorithm	Pagie-1	Kotanchek
Bingo	27%	22%	Bingo	<b>33</b> %	66%
EPLEX	28%	<b>18</b> %	EPLEX	45%	37%
GP-GOMEA	30%	<b>76</b> %	GP-	<b>100</b> %	<b>100</b> %
Operon	66%	74%	GOMEA		
PySR	36%	34%	Operon	<b>100</b> %	94%
SBP	42%	60%	PySR	52%	71%
			SBP	60%	<b>100</b> %



- Overparametrization is widespread in SymReg
- EqSat can consistently reduce the number of parameters
- · Sympy sometimes makes it worse
- Runtime is low most of the time and it can be improved

• Confidence Intervals with Profile Likelihood

