

Generating Implementations of Data Type Operations

Anthony Hunt and Emil Sekerinski

McMaster University,
Hamilton, Canada



Abstract Data Types

Operations [Abrial 10][Gries, Schneider 93]

Syntax	Label	Syntax	Label
$set(T)$	Unordered, unique collection	$S \cup T$	Union
$S \leftrightarrow T$	Relation, $set(S \times T)$	$S \cap T$	Intersection
$S \rightarrow \mathbb{N}$	Bag/Multiset	$S \setminus T$	Difference
\emptyset	Empty Set	$S \times T$	Cartesian Product
$\{x, y, \dots\}$	Set Enumeration	$dom(R)$	Domain
$\{x \cdot P \mid E\}$	Set Comprehension	$ran(R)$	Range
$x \mapsto y$	Pair	$R[S]$	Image
$S \nrightarrow T$	Partial Function	$R \boxplus Q$	Overriding
$S \rightarrowtail T$	Total Injective Function	$R \circ Q$	Composition
$S \subseteq T$	Subset	$S \triangleleft R$	Domain Restriction
$S - T$	Bag Difference	R^{-1}	Inverse

Visitor Information System

Motivating Examples

- Academics attend workshops throughout CAV 2025
- Every workshop must be held in one room
- Visitors may attend at most one workshop at a time

$Visitor = \{Shufang, Mark, Grigory\}$
 $Room = \{D346, D273, D305, D160, D152\}$
 $Workshop = \{SYNT, MOSCA, HCVS\}$

$location: Workshop \rightarrow Room$
 $attends: Visitor \rightarrow Workshop$

Total Injective Function

Partial Function

Then, the number of meals to prepare for a specific *room* would be:

$$card((location^{-1} \circ attends^{-1})[\{room\}])$$

Relational Image

Warehouse Inventory System

Motivating Examples

- Furniture material catalogue
- Warehouse inventory
- Product recipes

$Material = \{2 \times 4 \text{ Plank}, \text{Hex Bolt}, 3 \text{ Inch Screw}, \dots\}$

$Price = \mathbb{N}_0$

$Product = \{Cabinet, Desk, Bookshelf, \dots\}$

$catalogue: Material \rightarrow Price$

$inventory: bag[Material]$

$recipes: Product \rightarrow bag[Material]$

Restocking price, with a target inventory $inventory_t$:

$\sum p \mapsto n_m \cdot$

$p \mapsto n_m \in catalogue^{-1} \circ (inventory_t - inventory) \mid p \times n_m$

Bag Difference

Conway's Game of Life

Motivating Examples [Leuschel 20]

- Next state depends on previous position
- Cells with 2 or 3 neighbours remain living
- Dead cells with 3 living neighbours are revived

$Boundary \in \mathbb{N}$

$Living \in \mathbb{P}(Boundary \times Boundary)$



$neigh = \{x, y, x_n, y_n \cdot$

$$\begin{aligned} & x \in Boundary \wedge y \in Boundary \\ & \wedge x_n \in (x - 1)..(x + 1) \wedge y_n \in (y - 1)..(y + 1) \\ & \wedge x \neq x_n \wedge y \neq y_n \\ & | (x \mapsto y) \mapsto (x_n, \mapsto y_n) \} \end{aligned}$$

$nextStep =$

$$\begin{aligned} & \{ c \cdot c \in Living \wedge card(neigh[\{c\}] \cap Living) = 2 \mid c \} \\ \cup & \{ c \cdot c \in neigh[Living] \\ & \wedge card(neigh[\{c\}] \cap Living) = 3 \mid c \} \end{aligned}$$

Credit (Image): Conway's Game of Life. Wikipedia, 2025. https://en.wikipedia.org/wiki/Conway%27s_Game_of_Life

Abstract Data Types

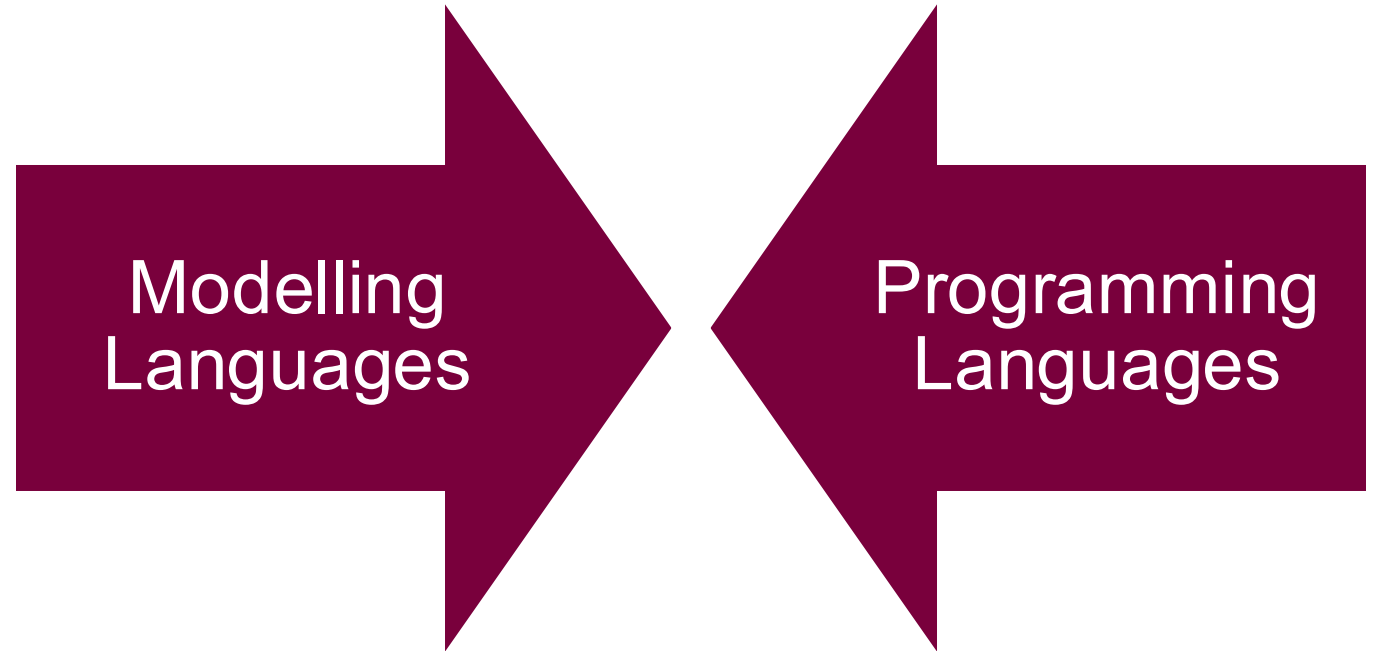
In Programming and Modelling Languages

		Sets						Relations											Efficiently Executable
		\cup	\cap	\setminus	\times	\subseteq	$\{x \cdot P \mid E\}$	dom ran	\cup	\cap	\setminus	\subseteq	$R[S]$	\boxtimes	\circ	\triangleleft	R^{-1}	Multiplicity	
Programming Languages	Python	✓	✓	✓	✓	✓	✓	✓	✗	*	*	*	*	✓	✗	✗	✗	✗	*
	Haskell	✓	✓	✓	✓	✓	*	✓	✗	✓	✓	✓	*	✓	*	*	✓	*	✓
	Rust	✓	✓	✓	*	✓	*	✓	✗	✗	✗	✗	✓	*	✗	✗	✗	✗	✓
	C	*	*	*	*	*	✗	*	*	*	*	*	*	*	*	*	*	✗	✓
	APL	✓	✓	✓	✓	*	✗	✓	*	*	*	*	*	*	✓	*	✓	✗	✓
Modelling Languages	SetL	✓	✓	✓	*	✓	✓	✓	✓	✓	✓	✓	✓	*	*	✓	*	*	*
	UML	✓	✓	✓	✓	✓	✗	✓	✗	✗	✗	✓	✓	✓	✗	✗	✗	✓	*
	Event-B	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	*
	Alloy	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✗

Goals

Of our set-based language

- Syntax close to discrete mathematics
- Simple, small, and usable
- Highly performant execution
- High-level set-theory optimization



Abstract Data Types in Python

Implementation of Set Operations in Modern Languages

$$(S \cup T) \cap V$$



$$(S \cap V) \cup ((T \setminus S) \cap V)$$

Python translation:

```
ret = set()
for x in S:
    ret.add(x)
# Step 1: len(ret) == len(S)
for x in T:
    ret.add(x)
# Step 2: len(ret) <= len(S) + len(T)
for x in ret:
    if x not in V:
        ret.remove(x)
# Step 3: len(ret) <= len(V)
```

New Python translation:

```
ret = set()
for x in S:
    if x in V:
        ret.add(x)
# Step 1: len(ret) <= min(len(S), len(V))
for x in T:
    if x not in S and x in V:
        ret.add(x)
# Step 2: len(ret) <= min(len(S) + len(T), len(V))
```


Related Work

- ProB animation engine for Event-B [Leuschel 22]
- GHC rewrite rules [Jones, Tolmach, Hoare 01]
- Rewrite systems for set-theory focused optimization [Cristia, Rossi 22][Vissier, Benaissa, Tolmach 98][Smith, Westfold 20]
- SETL data type *lowering* from abstract types to suitable concrete structures [Schonberg, Schwartz, Sharir 79][Freudenberger, Schwartz, Sharir 83]

Our Term Rewriting System

Overview: 7 Phases

- Set Comprehension Construction
- Predicate Disjunctive Normal Form
- Nesting
- Generator Selection
- Code Generation 1
- Equality Elimination
- Code Generation 2

Phase 1: Set Comprehension Construction

Set-typed variables and literals are decomposed into set comprehensions

Rewrite Rules	
Predicate Operations	$S \cup T \rightsquigarrow \{x \cdot x \in S \vee x \in T \mid x\}$
	$S \cap T \rightsquigarrow \{x \cdot x \in S \wedge x \in T \mid x\}$
	$S \setminus T \rightsquigarrow \{x \cdot x \in S \wedge x \notin T \mid x\}$
Membership Collapse	$x \in \{E \mid P\} \rightsquigarrow \exists y \cdot P \wedge x = E$
Image	$R[S] \rightsquigarrow \{y \cdot \exists x \cdot x \mapsto y \in R \wedge x \in S \mid y\}$
Composition	$x \mapsto z \in R \circ Q \rightsquigarrow \exists y, y' \cdot x \mapsto y \in R \wedge y' \mapsto z \in Q \wedge y = y'$
Inverse	$x \mapsto y \in R^{-1} \rightsquigarrow y \mapsto x \in R$
Cardinality	$\text{card}(S) \rightsquigarrow \sum x \cdot x \in S \mid 1$

Post-condition:

- All set-like terms must be comprehensions

location: Workshop \rightarrow Room
attends: Visitor \rightarrow Workshop

numMeals =
card((location⁻¹ \circ attends⁻¹)[{room}])



numMeals =
 $\sum v, r \cdot \exists p, p' \cdot v \mapsto p' \in \text{attends}$
 $\wedge p \mapsto r \in \text{location} \wedge p = p' \wedge r = \text{room} \mid 1$

Phase 2: Predicate Disjunctive Normal Form (DNF)

Predicate structure: top-level \vee -operations with nested \wedge -operations

Rewrite Rules

Flatten Nested \wedge	$x_1 \wedge \cdots \wedge (x_i \wedge x_{i+1}) \wedge \cdots$ $\rightsquigarrow x_1 \wedge \cdots \wedge x_i \wedge x_{i+1} \wedge \cdots$
Flatten Nested \vee	$x_1 \vee \cdots \vee (x_i \vee x_{i+1}) \vee \cdots$ $\rightsquigarrow x_1 \vee \cdots \vee x_i \vee x_{i+1} \vee \cdots$
Distribute \wedge over \vee	$x \wedge (y \vee z) \rightsquigarrow (x \wedge y) \vee (x \wedge z)$
Double Negation	$\neg \neg x \rightsquigarrow x$
De Morgan	$\neg(x \wedge y) \rightsquigarrow \neg x \vee \neg y$
	$\neg(x \vee y) \rightsquigarrow \neg x \wedge \neg y$

Post-condition:

- All set-like terms must be comprehensions
- Quantifier predicates are in DNF

$numMeals =$

$$\sum v, r \cdot \exists p, p' \cdot v \mapsto p' \in attends$$
$$\wedge p \mapsto r \in location \wedge p = p' \wedge r = room \mid 1$$

Phase 3: Nesting

Restructuring quantifications to better suit code generation

Rewrite Rules

Nesting

$\{x, y \cdot P \wedge Q \mid E\} \rightsquigarrow \{x \cdot P \mid \{y \cdot Q \mid E\}\}$
Where y does not occur in P

- Nesting selection is currently arbitrary; in the future we plan to use better selection heuristics

Post-condition:

- All set-like terms must be comprehensions
- Quantifier predicates are in DNF
- One bound variable per quantifier

$numMeals =$

$\sum v, r \cdot \exists p, p' \cdot v \mapsto p' \in attends$
 $\wedge p \mapsto r \in location \wedge p = p' \wedge r = room \mid 1$



$numMeals =$

$\sum r \cdot \exists p \cdot p \mapsto r \in location \wedge r = room$
 $\mid (\sum v \cdot \exists p' \cdot v \mapsto p' \in attends \wedge p = p' \mid 1)$

Phase 4: Generator Selection

Selecting element generators for use as iterables in generated for-loops

Rewrite Rules

Generator Selection

$$\bigwedge P_i \rightsquigarrow P_g \wedge \bigwedge_{P_i \neq P_g} P_i$$

Where:

- P_g is of the form $x \in S$
- x is the quantifier's bound variable
- S is set-like term

Post-condition:

- All set-like terms must be comprehensions
- Quantifier predicates are in DNF
 - Guaranteed top-level-v operation
- One bound variable per quantifier
- All predicate v-clauses have an assigned generator

- In the case of multiple candidate generators, selection is arbitrary
- Generator Selection is only valid inside a quantifier's predicate

numMeals =

$$\sum r \cdot \exists p \cdot p \mapsto r \in location \wedge r = room \\ | (\sum v \cdot \exists p' \cdot v \mapsto p' \in attends \wedge p = p' \mid 1)$$



numMeals =

$$\sum \mathbf{r} \cdot \exists p \cdot \mathbf{p} \mapsto \mathbf{r} \in \mathbf{location} \wedge r = room \\ | (\sum \mathbf{v} \cdot \exists p' \cdot \mathbf{v} \mapsto \mathbf{p}' \in \mathbf{attends} \wedge p = p' \mid 1)$$

Phase 5: Code Generation 1

Start lowering expressions into imperative-like loops

Rewrite Rules

Quantifier Generation	$\oplus E \mid P$	$a := \text{identity}(\oplus)$ $\rightsquigarrow \text{loop } P \text{ do}$ $a := a \oplus E$
Disjunct Conditional	$\text{loop } \bigvee P_i \text{ do}$ $body$	$\rightsquigarrow \text{loop } P_0 \text{ do}$ $body$ $\rightsquigarrow \text{loop } P_1 \wedge \neg P_0 \text{ do}$ $body$ \dots
Nested Quantifier	$a := a \oplus (\oplus E \mid P)$	$\rightsquigarrow \text{loop } P \text{ do}$ $a := a \oplus E$

Post-condition:

- No quantifiers exist within the AST
- No v-operators exist within a *loop*'s predicate
- All predicate v-clauses have an assigned generator

- \oplus is any of $\{\dots\}, \Sigma, \Pi, \cup, \cap$

$$\begin{aligned}
 \text{numMeals} = & \\
 & \Sigma \mathbf{r} \cdot \exists p \cdot \mathbf{p} \mapsto \mathbf{r} \in \text{location} \wedge r = \text{room} \\
 & \mid (\Sigma \mathbf{v} \cdot \exists p' \cdot \mathbf{v} \mapsto \mathbf{p}' \in \text{attends} \wedge p = p' \mid 1)
 \end{aligned}$$



```

numMeals := 0
loop p  $\mapsto$  r  $\in$  location  $\wedge$   $r = \text{room}$  do
  loop v  $\mapsto$  p'  $\in$  attends  $\wedge$   $p = p'$  do
    numMeals := numMeals + 1
  
```

Phase 6: Equality Elimination

Eliminate all undefined variables (with implicit \exists quantifiers)

Rewrite Rules

Equality
Elimination

$$\text{loop } \bigwedge P_i \text{ do } \text{body} \rightsquigarrow \text{loop } \bigwedge_{P_i \neq P_{eq}} P_i[x := E] \text{ do } \text{body}[x := E]$$

Where:

- P_{eq} is of the form $x = E$
- x is undefined and unbound in the current scope

Post-condition:

- No quantifiers exist within the AST
- No v-operators exist within a *loop*'s predicate
- All predicate v-clauses have an assigned generator
- All *loop* predicate variables are defined.

```
numMeals := 0
loop p  $\mapsto$  r  $\in$  location  $\wedge$   $r = \text{room}$  do
  loop v  $\mapsto$  p'  $\in$  attends  $\wedge$   $p = p'$  do
    numMeals := numMeals + 1
```

Phase 7: Code Generation 2

Eliminate all intermediate *loop* structures

Rewrite Rules

Conjunct
Conditional

$$\begin{array}{l} \text{loop } P_g \\ \wedge \bigwedge P_i \text{ do } \rightsquigarrow \\ \text{body} \end{array} \quad \rightsquigarrow \quad \begin{array}{l} \text{if } \bigwedge_{\text{free}(P_i)} P_i \text{ then} \\ \text{for } P_g \text{ do} \\ \text{if } \bigwedge_{\neg \text{free}(P_i)} P_i \text{ then} \\ \text{body} \end{array}$$

```
numMeals := 0
loop p ↦ r ∈ location ∧ r = room do
  loop v ↦ p' ∈ attends ∧ p = p' do
    numMeals := numMeals + 1
```



```
numMeals := 0
for p, r ∈ location do
  if r = room then
    for v, p' ∈ attends do
      if p = p' then
        numMeals := numMeals + 1
```

Post-condition:


- Code is imperatively executable

Visitor Information System - Translation


Motivating Examples

$location: Workshop \mapsto Room$
 $attends: Visitor \nrightarrow Workshop$

$$\begin{aligned} numMeals \\ &= card((location^{-1} \circ attends^{-1})[\{room\}]) \end{aligned}$$


$$\begin{aligned} numMeals \\ &= \sum v, r \cdot \exists p, p' \cdot v \mapsto p' \in attends \\ &\quad \wedge p \mapsto r \in location \wedge \mathbf{p} = \mathbf{p}' \wedge \mathbf{r} = \mathbf{room} \mid 1 \end{aligned}$$


$$\begin{aligned} numMeals \\ &= \sum v \cdot \exists p' \cdot v \mapsto p' \in attends \wedge location(p') = room \mid 1 \end{aligned}$$



```
numMeals := 0
for v, p' in attends do
  if location(p') = room then
    numMeals := numMeals + 1
```


Warehouse Inventory System – Translation

Motivating Examples

catalogue: *Material* \rightarrow *Price*
inventory: *bag*[*Material*]
recipes: *Product* \rightarrow *bag*[*Material*]

Restocking price, with a target inventory *inventory_t*:



$price = \sum p \mapsto n_m \cdot$
 $p \mapsto n_m \in catalogue^{-1} \circ (inventory_t - inventory) \mid p \times n_m$

$price = \sum m \mapsto n_m \cdot \exists n, m, p \cdot m \mapsto n_m \in inventory_t$
 $\wedge n = n_m - inventory(m) \wedge n \geq 0 \wedge p = catalogue(m) \mid p \times n_m$

price := 0
for *m*, *n_m* \in *inventory_t* **do**
 n := *n_m* - *inventory*(*m*)
 if *n* \geq 0 **then**
 price := *price* + *catalogue*(*m*) \times *n*

Conway's Game of Life – Translation

Motivating Examples

$$\begin{aligned} \text{neigh} = & \{x, y, x_n, y_n \cdot \\ & x \in \text{Boundary} \wedge y \in \text{Boundary} \\ & \wedge x_n \in (x - 1) .. (x + 1) \\ & \wedge y_n \in (y - 1) .. (y + 1) \\ & \wedge x \neq x_n \wedge y \neq y_n \\ & \mid (x, y) \mapsto (x_n, y_n)\} \end{aligned}$$
$$\begin{aligned} \text{nextStep} = & \\ & \{c \cdot c \in \text{Living} \wedge \text{card}(\text{neigh}[\{c\}] \cap \text{Living}) = 2 \mid c\} \\ & \cup \{c \cdot c \in \text{neigh}[\text{Living}] \\ & \quad \wedge \text{card}(\text{neigh}[\{c\}] \cap \text{Living}) = 3 \mid c\} \end{aligned}$$

$$\begin{aligned} \text{nextStep} = & \\ & \{c \cdot c \in \text{Boundary} \times \text{Boundary} \\ & \wedge \exists n \cdot n = \text{card}(\text{neigh}[\{c\}] \cap \text{Living}) \\ & \wedge ((c \in \text{Living} \wedge n = 2) \vee n = 3) \mid c\} \end{aligned}$$
$$\begin{aligned} \text{nextStep} & := \emptyset \\ \text{for } i \in \text{Boundary} & \text{ do} \\ & \text{for } j \in \text{Boundary} \text{ do} \\ & \quad n := 0 \\ & \quad \text{for } n' \in \text{neigh}[\{i \mapsto j\}] \text{ do} \\ & \quad \quad \text{if } n' \in \text{Living} \text{ then} \\ & \quad \quad \quad n := n + 1 \\ & \quad \text{if } n = 3 \vee (n = 2 \wedge i \mapsto j \in \text{Living}) \text{ then} \\ & \quad \quad \text{nextStep} := \text{nextStep} \cup \{i \mapsto j\} \end{aligned}$$


Preliminary Results

- Intermediate sets never grow larger than $\max(\text{starting sets}, \text{resulting set})$
- Unions, Relation Overriding, and Cartesian Products are limited in growth
- Other operators only decrease the size of a resulting set
- What about running time?...

Preliminary Results

Operation Running Time

Operation	Expected Running Time
$S \cup T$	$O(S + T)$
$S \cap T$	$O(\min(S , T))$
$S \times T$	$O(S T)$
$R[S]$	$O(\min(S , R))$
$S \triangleleft R$	$O(R)$
$R \boxtimes Q$	$O(\max(Q , R))$
$R \circ Q$	$O(R Q)$

Preliminary Results

Examples Running Time

	Core Expression	Naïve Running Time	Optimized Running Time
Visitor Information System	$(location^{-1} \circ attends^{-1})[\{room\}]$	$O(location attends)$	$O(attends)$
Warehouse Inventory System	$\sum p \mapsto n_m \cdot p \mapsto n_m$ $\in (inventory_t - inventory) \circ catalogue^{-1}$ $ p \times n_m$	$O(inventory_t catalogue)$	$O(inventory_t)$
Conway's Game of Life	$\{c \cdot c \in Living \text{ —————}$ $\wedge card(neigh[\{c\}] \cap Living) = 2 \mid c\}$ $\cup \{c \cdot c \in neigh [Living]$ $\wedge card(neigh [\{c\}] \cap Living) = 3 \mid c\}$	$O(2 Boundary ^2)$	$O(Boundary ^2)$

$Living \in \mathbb{P}(Boundary \times Boundary)$

Preliminary Results

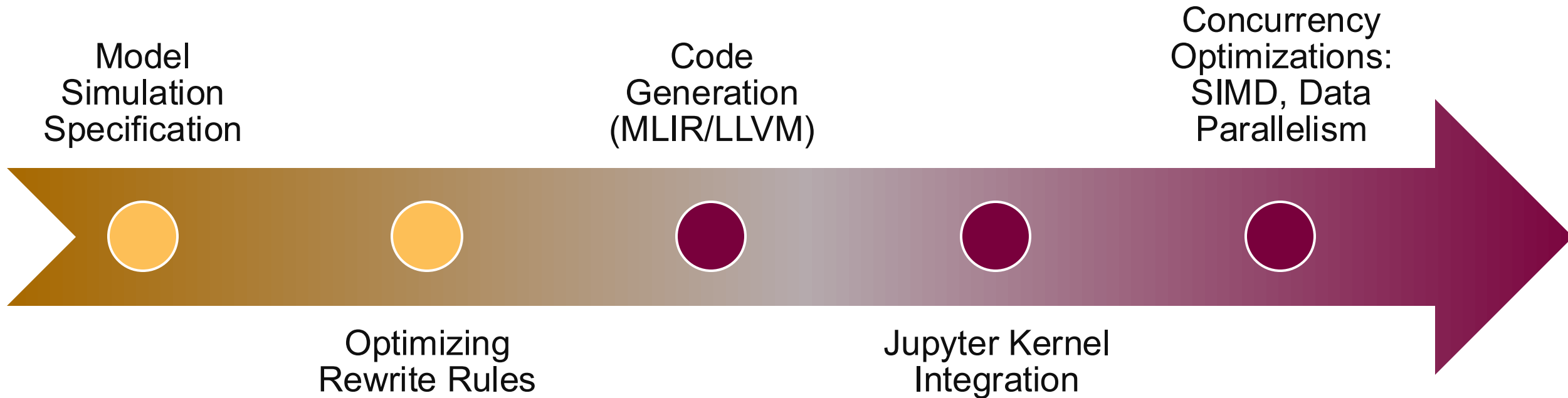
Complex Operation Running Time

Operation	Expected Running Time
$(S \cup T) \cap U$	$O(\min(S + T , U))$
$S \cap R[T]$	$O(\min(S , R , T))$
$(R \circ Q)[S]$	$O(\min(S , R) + \min(R[S] , Q))$

- Simplify large groups of \wedge -predicates
- Eagerly apply conditions on sequential \vee -predicates

Current Compiler State and Roadmap

Development of a High-Level, Efficient, Set-Based Language



References

- Michael Leuschel. *ProB2 Jupyter Notebooks* *GitLab*. 2020. https://gitlab.cs.uni-duesseldorf.de/general/stups/prob2-jupyter-notebooks/-/blob/9d830112d9713d0cbc12d9d285f44ee61c827ccf/puzzles/Game_of_Life.ipynb
- Jean-Raymond Abrial. *Modeling in Event-B: system and software engineering*. Cambridge University Press, 2010.
- David Gries and Fred Schneider. *A Logical Approach to Discrete Math*. Springer New York, 1993.
- Robert Dewar. *The SetL Programming Language*. Bell Laboratories, 1979.
- Simon Peyton Jones, Andrew Tolmach, and Tony Hoare. *Playing by the rules: rewriting as a practical optimisation technique in GHC*. In: 2001 Haskell Workshop. ACM SIGPLAN. Sept. 2001.
- Michael Leuschel. *Programming in B: Sets and Logic all the Way Down*. In: LPOP 2022.
- Maximiliano Cristia and Gianfranco Rossi. *{log}: Programming and Automated Proof in Set Theory*. In: LPOP 2022.
- Douglas R. Smith and Stephen J. Westfold. *Transformations for Generating Type Refinements*. In: Formal Methods. FM 2019 International Workshops. Springer International Publishing, 2020, pp. 371-387
- Edmond Schonberg, Jacob T. Schwartz, and Micha Sharir. 1979. *Automatic data structure selection in SETL*. In Proceedings of the 6th ACM SIGACT-SIGPLAN symposium on Principles of programming languages (POPL '79). Association for Computing Machinery, New York, NY, USA, 197–210. <https://doi.org/10.1145/567752.567771>
- Stefan M. Freudenberger, Jacob T. Schwartz, and Micha Sharir. 1983. *Experience with the SETL Optimizer*. ACM Trans. Program. Lang. Syst. 5, 1 (Jan. 1983), 26–45. <https://doi.org/10.1145/357195.357197>