

Types in Programming Languages

Alexey Chentsov

Contents

Introduction	3
1 Simple systems	8
1.1 Untyped arithmetics	8
1.1.1 Language and inductive definitions	8
1.1.2 Induction on terms	10
1.1.3 Semantic styles	12
1.1.4 Evaluation	13
1.2 Typed arithmetic	17
1.2.1 Typing relation	17
1.2.2 Type safety	19
1.3 Untyped λ -calculus	20
1.3.1 Lambda terms	20
1.3.2 Beta reduction	22
1.3.3 Programming in pure λ -calculus	23
1.3.4 Formal definitions	24
1.4 Simply-typed λ -calculus	24
1.4.1 Properties of λ_{\rightarrow}	24
1.4.2 Type reconstruction problem	26
2 Polymorphism	29
2.1 Introduction	29
2.1.1 Ad hoc polymorphism	30
2.1.2 Subtype polymorphism	31
2.1.3 Introduction to parametric polymorphism	33
2.2 Let-polymorphism	35
2.3 First-class polymorphism	39
3 Selected type-based programming techniques	43
3.1 Elements of generic programming	43
3.1.1 Basic notions of generic programming	44

3.1.2	Example of algorithm generalization	44
3.2	Datatype-generic programming	49
3.2.1	Roy–Floyd–Warshall algorithm	50
3.2.2	Algebras & Morphisms	52
3.2.3	Type constructors and functors	55
3.2.4	Polynomial Functors	56
3.2.5	Functor algebras	59
	Bibliography	62
	Glossary	64

Introduction

Let us start by considering several systems.

1. Bit strings: computer data (in general sense).
2. S-expressions: recursive structure loosely based on lists (ordered pairs and atoms)
3. λ -expressions.
4. Sets: unstructured collections of abstract elements.

They can be thought as universes. For flexibility, universality and simplicity they are untyped or actually of one type. For instance, the main difference between the objects in the first case is the length of the bit string.

On the other hand from practical standpoint it is convenient to distinguish different sorts of objects due to manipulations that are allowed with them or their behavior. For instance, operations over integers and doubles are performed by different units. Then there should be a *representation*¹ of these entities in universal form. But different kinds of objects might have the same representations. Examples:

1. Integers, floating point numbers, strings, instructions (code) look the same.
2. Both program code and data (in lisp) is represented as S-expression.
3. In pure λ -expressions everything is a function.
4. Sets are a universal mathematical language that allows building very complex mathematical structures.

¹Ukrainian equivalent: *представлення*

Types are not only the matter of convenience, for naive (untyped) set theory a number of paradoxes, e.g. Russel paradox, arise and one of the ways to avoid them is to introduce two types: proper classes and sets.

So in general the type system is used to organize and systematize the object universe.

For typed universe many constructions lack the meaning. For instance, very few bit strings can represent a meaningful program, neither does S-expression (1 'three +), concatenation of word 'one' and negative number -0.5 doesn't make sense as well.

Type systems make the programmer more aware of the conventional rules regarding valid manipulation with different entities.

In the case of programming languages type system is used to avoid constructions that don't have meaning. As Luca Cardelli puts it[1]

Type is a suit of armor that protects the representation from arbitrary or unintended use.

Basically, using types is similar to using units of measurements in physics. Only commensurable quantities might be added, subtracted, compared. It can be used as sanity check for expressions (dimensional homogeneity combined with factor-label method).

There is a loose definition of *type system* [2, p. 1]

A *type system* is a tractable syntactic method for proving the *absence* of certain *program* behaviors by *classifying* phrases according to the kinds of *values* they compute.

Type systems are popular and most established lightweight formal methods.

Here we come to instrumental essence of the type system. It allows us to avoid special kind of errors *runtime type errors*.

“Classifying according to kind of values” means that type system may be regarded as calculating approximation to the runtime behaviour. In this respect *type checking*² is *conservative*. Type system makes sure that runtime type errors are *absent* but it cannot check if they are *present*. E.g.

if <complex test> then 5 else <type error>

Even though the condition might always be true and else branch would be never executed the program is rejected by type system (due to limitations). This situation is similar to banker's algorithm when state is safe only when absence of deadlock could be guaranteed.

²Ukrainian equivalent: *перевірка мунію*

Dynamic vs static type checking We can talk about *manifest* and *latent* properties or attributes. Attributes known at compile time are called manifest and known at run-time are called latent. For instance, the Lvalue (i.e. the variable) in, say, Java has certain known type T and its Rvalue (i.e. content) has unknown type which is a subtype of T .

Remark. The distinction between manifest and latent attributes is not clear cut and subject to change with evolution of hardware and programming languages. For instance, in C++ expression that previously were determined on execution now can be computed by compile time functions (const-expressions).

Based on this we can do *static* type checking³ and *dynamic* type checking⁴. Type requirements explicitly placed in the code by programmer are called *type annotations*. These could vary from quite modest (Haskell, ML) to very redundant (C, Pascal). In the case of very extensive type annotations type checking effectively turns into proof checking.

Dynamic type checking⁵ requires metainformation about typesystem in runtime (RTTI).

Remark. Luca Cardelli distinguishes *explicitly typed*, *implicitly typed* and *untyped* languages[3]. Type system is the static type system. Untyped languages provide safety with *dynamic checks*. Typed languages use both static tests and dynamic tests (based on the static type system) to provide safety.

The whole type checking can be performed at runtime. However modern programming languages prefer static type checking for two reasons:

- early “availability”: mistakes are found before program execution;
- more effective code: dynamic type checking involves extra work in runtime.

On the other hand, static type checking might interfere with sound programming techniques that are incompatible with early binding [1].

Remark. Actually, if CPU were equipped with type determination facilities the performance would have been no longer an issue for dynamic type checking.

³Ukrainian equivalent: *статична* перевірка типів

⁴Ukrainian equivalent: *динамічна* перевірка типів

⁵sometimes called *dynamic typing*, arguably a misnomer

Language safety The main motif for type system as formulated by Robin Milner is “Well-typed programs cannot go wrong”. It is closely related with notion of language safety. This term is even more contentious than type system. There are several definitions of *language safety*. According to Benjamin Pierce

A safe language is one that protects its own abstractions [2].

This means it is not necessary to keep in mind all sorts of low-level implementation details, because these abstractions are consistently and completely defined in terms of their own properties [4].

Cardelli distinguished two types of errors. The ones that are detected right away on the spot are called *trapped* and the ones that go unnoticed and then lead to arbitrary behavior called *untrapped*. An example of the latter is stack overflow. Then

A safe language prevents untrapped run-time errors [3].

Finally, there is a definition

A safe language is completely defined by its programmer’s manual [2].

It means that behavior of any program can be predicted fully based on the language manual. For instance, correctness of some programming techniques of unsafe language might depend on calling convention or stack layout details.

Functions of type system The most important advantage of having a typesystem is code maintainability [3].

Besides error detection, additionally type systems (type annotations) serve other purposes:

- abstraction
- documentation
- efficiency.

It provides better interfaces between the modules (partial contract) which form the basis for more *abstract* design. It *documents* code and it is always up-to-date due to check on every compilation. Statically checked programs have better *efficiency* because of the absence of dynamic checks and optimization based on properties of the types. The latter is mostly related to base types: hard-coding floating point operations, loop unrolling etc. Today, high-performance compilers heavily rely on information gathered by type-checker. One of the strategic goals in this direction is to allow optimizations for user-defined types similar to base types.

Note on supplementary code The textbook [2] is accompanied with implementations of all discussed systems (languages) starting with untyped boolean arithmetic and finishing with system F extended with subtyping. Supplementary systems are provided in OCaml programming language. They can be considered as (formal) language processors. Therefore parser is one of their key components. Choice of programming language is due to author's personal preferences but also due to native support of inductive types in OCaml which simplifies construction and representation of *abstract syntax tree*⁶.

Provided systems are designed to operation in batch mode. This means they expect input in the form of user provided file. This input contains expressions ("lines of code") of the given systems. The system then forms the result of expression evaluation on the standard output. For some tasks batch mode is not the most convenient. To experiment with the system the interactive mode is preferred. Also access to the source code allows system extension. Some system are already enriched with the elements of the others. For instance, basic arithmetic is embedded in the most of the system. These extensions implies some tweaks to the parser and maybe evaluation component of the system. Particularly it is appropriate for the varieties of lambda calculus implementation to use interactive mode and add syntactic sugar to allow conventional function definition syntax. This is utilized in this particular course.

⁶Ukrainian equivalent: *абстрактне синтаксичне дерево*

Chapter 1

Simple systems

In order to discuss questions concerning the types in (modern) programming languages one has to present the formal description of systems used. Since immediate exploration of advanced programming constructions could be too much of a leap in complexity we would start from simple systems and gradually move to more modern counterparts. It is important to mention that one of the systems provides the basis for many programming languages and thus plays very important role in theoretical computer science.

1.1 Untyped arithmetics

1.1.1 Language and inductive definitions

Description of the (computing) system begins with the formal description of the language. We have several ways to do this. Most conventional is to use grammars and Backus–Naur form (or its variation).

<code>t ::=</code>		<i>terms:</i>
<code> true</code>		<i>constant true</i>
<code> false</code>		<i>constant false</i>
<code> if t then t else t</code>		<i>conditional</i>
<code> 0</code>		<i>constant zero</i>
<code> succ t</code>		<i>successor</i>
<code> pred t</code>		<i>predecessor</i>
<code> iszero t</code>		<i>zero test</i>

In the l.h.s. we specify that this defines the *terms* of the language. In the r.h.s. we have *alternatives*. In this construction `t` in the r.h.s. can be

substituted with any correct term of the language. Since it is not the part of the language it is a *meta*-variable. We see that more complex terms are based on simple terms (*subterms*). This is an example of inductive definition. In mathematics this construction is introduced using *inductively defined sets*. Here is the definition

1. $\{\text{true}, \text{false}, 0\} \subseteq \mathcal{T}$;
2. if $t_1 \in \mathcal{T}$, then $\{\text{succ } t_1, \text{pred } t_1, \text{iszero } t_1\} \subseteq \mathcal{T}$;
3. if $t_1 \in \mathcal{T}$, $t_2 \in \mathcal{T}$, $t_3 \in \mathcal{T}$, then $\text{if } t_1 \text{ then } t_2 \text{ else } t_3 \in \mathcal{T}$.

First point gives *base clause* (*basis step*). The second and third points are *inductive clause* (*recursive step*). Inductive definition also has implicit part known as *extremal clause*¹ which states that “only elements satisfying base and inductive clause are in the set” or “nothing else is an element of the set”. It guarantees uniqueness of the set and thus the correctness of the definition. Extremal clause is usually omitted.

Preferable way to formulate this construction in theoretical computer science is based on *inference rules*. Inference rule states that if premises are true then the conclusion is too. Inference rule are written down as fractions where numerator is premises and denominator is conclusion of the rule. The rules with empty premises are axioms. This notation is often used in natural deduction systems². The previous definition would take the form

$$\begin{array}{ccc}
 \text{true} \in \mathcal{T} & \text{false} \in \mathcal{T} & 0 \in \mathcal{T} \\
 \frac{t_1 \in \mathcal{T}}{\text{succ } t_1 \in \mathcal{T}} & \frac{t_1 \in \mathcal{T}}{\text{pred } t_1 \in \mathcal{T}} & \frac{t_1 \in \mathcal{T}}{\text{iszero } t_1 \in \mathcal{T}} \\
 \\
 \frac{t_1 \in \mathcal{T} \quad t_2 \in \mathcal{T} \quad t_3 \in \mathcal{T}}{\text{if } t_1 \text{ then } t_2 \text{ else } t_3 \in \mathcal{T}}
 \end{array}$$

Finally, there is a non-inductive way to define the terms.

$$\begin{aligned}
 S_0 &= \emptyset \\
 S_{i+1} &= \{\text{true}, \text{false}, 0\} \\
 &\quad \cup \{\text{succ } t_1, \text{pred } t_1, \text{iszero } t_1 \mid t_1 \in S_i\} \\
 &\quad \cup \{\text{if } t_1 \text{ then } t_2 \text{ else } t_3 \mid t_1, t_2, t_3 \in S_i\}. \\
 \\
 S &= \bigcup_{i \in \mathbb{N}} S_i.
 \end{aligned}$$

¹rarely, *closure clause*

²The notation is attributed to Gentzen

Proposition 1.1. *Family $\{S_i\}_{i \in N}$ form a cumulative hierarchy, that is*

$$\forall i \in N. S_i \subseteq S_{i+1}.$$

Proposition 1.2. $\mathcal{T} = S$.

Proof. It is enough to prove the characteristic properties of \mathcal{T} for S . First, let us check that conditions 1–3 hold for S . By definition $\{\mathbf{true}, \mathbf{false}, 0\} = S_1 \subseteq S$, i.e. 1 holds. Suppose $\mathbf{t} \in S$. Then

$$\mathbf{t} \in S_i \Rightarrow \mathbf{succ} \mathbf{t}, \mathbf{pred} \mathbf{t}, \mathbf{iszero} \mathbf{t} \in S_{i+1} \Rightarrow \mathbf{succ} \mathbf{t}, \mathbf{pred} \mathbf{t}, \mathbf{iszero} \mathbf{t} \in S.$$

Finally if $\mathbf{t}_{1,2,3} \in S$ then

$$\begin{aligned} \mathbf{t}_k \in S_{i_k} \Rightarrow \mathbf{t}_k \in S_m \wedge m = \max\{i_1, i_2, i_3\} \Rightarrow \mathbf{if} \mathbf{t}_1 \mathbf{then} \mathbf{t}_2 \mathbf{else} \mathbf{t}_3 \in S_{m+1} \\ \Rightarrow \mathbf{if} \mathbf{t}_1 \mathbf{then} \mathbf{t}_2 \mathbf{else} \mathbf{t}_3 \in S. \end{aligned}$$

Now suppose that S' satisfies conditions 1–3 (of \mathcal{T}). Let us prove that $S \subseteq S'$. It is enough to show that $S_i \subseteq S'$ for all i . We will show that by induction. Obviously $S_0 = \emptyset \subseteq S'$. Let $\mathbf{t} \in S_{i+1}$ for some i . We have to consider several (general) cases. Case 1: $\mathbf{t} \in S_1$. Then $\mathbf{t} \in S$ by condition 1. Case 2: $\mathbf{t} = \mathbf{pred} \mathbf{s}$, $\mathbf{s} \in S_i$. Then by inductive hypothesis $\mathbf{s} \in S'$. As a consequence $\mathbf{t} \in S'$ by property 2. Case 3: $\mathbf{t} = \mathbf{if} \mathbf{t}_1 \mathbf{then} \mathbf{t}_2 \mathbf{else} \mathbf{t}_3$, $\mathbf{t}_i \in S_i$. Then by inductive hypothesis $\mathbf{t}_i \in S'$. As a consequence $\mathbf{t} \in S'$ by property 3. \square

This proposition reaffirms that terms of the language \mathcal{T} have limited number of forms: term is either (1) a constant, (2) has the form $\mathbf{succ} \mathbf{t}$, $\mathbf{pred} \mathbf{t}$, $\mathbf{iszero} \mathbf{t}$, or (3) a conditional $\mathbf{if} \mathbf{t}_1 \mathbf{then} \mathbf{t}_2 \mathbf{else} \mathbf{t}_3$.

1.1.2 Induction on terms

The characterization of terms of \mathcal{T} allows special definitions and reasoning on them. First, we introduce *inductive definitions* of functions. Let us define three functions. $\mathit{consts}(\mathbf{t})$ is a multiset of constants used in \mathbf{t} :

$$\begin{aligned} \mathit{consts}(\mathbf{true}) &= \{\mathbf{true}\} \\ \mathit{consts}(\mathbf{false}) &= \{\mathbf{false}\} \\ \mathit{consts}(0) &= \{0\} \\ \mathit{consts}(\mathbf{succ} \mathbf{t}_1) &= \\ &\quad \mathit{consts}(\mathbf{pred} \mathbf{t}_1) \\ &\quad \mathit{consts}(\mathbf{iszero} \mathbf{t}_1) &= \mathit{consts}(\mathbf{t}_1) \\ \mathit{consts}(\mathbf{if} \mathbf{t}_1 \mathbf{then} \mathbf{t}_2 \mathbf{else} \mathbf{t}_3) &= \mathit{consts}(\mathbf{t}_1) \cup \mathit{consts}(\mathbf{t}_2) \cup \mathit{consts}(\mathbf{t}_3) \end{aligned}$$

where operation \cup in the last line is union of multisets. Size of the term (number of node in abstract syntax tree of \mathfrak{t}):

$$\begin{aligned}
size(\mathfrak{true}) &= \\
size(\mathfrak{false}) &= \\
size(0) &= 1 \\
size(\mathfrak{succ } \mathfrak{t}_1) &= \\
size(\mathfrak{pred } \mathfrak{t}_1) &= \\
size(\mathfrak{iszero } \mathfrak{t}_1) &= size(\mathfrak{t}_1) + 1 \\
size(\mathfrak{if } \mathfrak{t}_1 \mathfrak{ then } \mathfrak{t}_2 \mathfrak{ else } \mathfrak{t}_3) &= size(\mathfrak{t}_1) + size(\mathfrak{t}_2) + size(\mathfrak{t}_3) + 1
\end{aligned}$$

Depth of the term \mathfrak{t} :

$$\begin{aligned}
depth(\mathfrak{true}) &= \\
depth(\mathfrak{false}) &= \\
depth(0) &= 1 \\
depth(\mathfrak{succ } \mathfrak{t}_1) &= \\
depth(\mathfrak{pred } \mathfrak{t}_1) &= \\
depth(\mathfrak{iszero } \mathfrak{t}_1) &= depth(\mathfrak{t}_1) + 1 \\
depth(\mathfrak{if } \mathfrak{t}_1 \mathfrak{ then } \mathfrak{t}_2 \mathfrak{ else } \mathfrak{t}_3) &= \max(depth(\mathfrak{t}_1), depth(\mathfrak{t}_2), depth(\mathfrak{t}_3)) + 1
\end{aligned}$$

Second, let us consider new inductive methods of reason that suits our inductively defined sets.

Theorem 1.3 (Induction on terms). *Suppose P is a predicate on terms.*

Induction on depth:

If, for each s ,
given $P(\mathfrak{r})$ for all \mathfrak{r} such that $depth(\mathfrak{r}) < depth(s)$
we can show $P(s)$,
then $P(s)$ holds for all s .

Induction on size:

If, for each s ,
given $P(\mathfrak{r})$ for all \mathfrak{r} such that $size(\mathfrak{r}) < size(s)$ we can show $P(s)$,
then $P(s)$ holds for all s .

Structural induction:

If, for each s ,
given $P(\mathfrak{r})$ for all immediate subterms \mathfrak{r} of s we can show $P(s)$,
then $P(s)$ holds for all s .

Structural induction was introduced by Burstall (1969) [5]. This technique was further popularized in functional programming. We demonstrate its application on the following

Proposition 1.4. $|const(\mathfrak{t})| \leq size(\mathfrak{t})$.

Proof. By structural induction (by induction on \mathfrak{t}). We have to consider three cases.

Case: $\mathfrak{t} \in \{\mathbf{true}, \mathbf{false}, 0\}$

By definition $|const(\mathfrak{t})| = 1 \leq 1 = size(\mathfrak{t})$.

Case: $\mathfrak{t} = \mathbf{succ} \mathfrak{t}_1, \mathbf{pred} \mathfrak{t}_1$, or $\mathbf{iszero} \mathfrak{t}_1$

Then $|const(\mathbf{succ} \mathfrak{t}_1)| = |const(\mathfrak{t}_1)| \leq size(\mathfrak{t}_1) < size(\mathfrak{t}_1) + 1 = size(\mathfrak{t})$, where induction hypothesis for \mathfrak{t}_1 is used.

Case: $\mathfrak{t} = \mathbf{if} \mathfrak{t}_1 \mathbf{then} \mathfrak{t}_2 \mathbf{else} \mathfrak{t}_3$

Then

$$\begin{aligned} |const(\mathfrak{t})| &= |const(\mathfrak{t}_1) \cup const(\mathfrak{t}_2) \cup const(\mathfrak{t}_3)| \\ &\leq |const(\mathfrak{t}_1)| + |const(\mathfrak{t}_2)| + |const(\mathfrak{t}_3)| \\ &\leq size(\mathfrak{t}_1) + size(\mathfrak{t}_2) + size(\mathfrak{t}_3) \\ &< size(\mathfrak{t}) \end{aligned}$$

where inductive hypothesis for $\mathfrak{t}_{1,2,3}$ is used. □

So in general technique boils down to proving the property for term based on inductive hypothesis that it holds for “smaller” terms. This requires considering a number of base and inductive cases.

1.1.3 Semantic styles

So far we have only talked about structure of the language expressions. Now we should provide (computational) meaning of the expressions, i.e. language semantics. There are three main approaches to define semantics called *semantic styles*³.

Axiomatic semantics can be traced back to the work of Floyd and later Hoare. The meaning of the program is defined in terms of assertions and their relationships. There are (logical) rules assigned to each base constructions of the language that specify their effect on pre- and post-conditions. Combining these rules one can reason about the effect of the whole program. Thus the meaning of programs is just what we can prove about them.

³Ukrainian equivalent: *семантичні стилі*

B (*untyped*)

Syntax		Evaluation	
$t ::=$	<code>true</code> <code>false</code> <code>if t then t else t</code>	<i>terms:</i> constant true constant false conditional	$t \longrightarrow t'$ <code>if true then t₂ else t₃ \longrightarrow t₂</code> (E-IFTRUE) <code>if false then t₂ else t₃ \longrightarrow t₃</code> (E-IFFALSE) $\frac{t_1 \longrightarrow t'_1}{\text{if } t_1 \text{ then } t_2 \text{ else } t_3 \longrightarrow \text{if } t'_1 \text{ then } t_2 \text{ else } t_3}$ (E-IF)
$v ::=$	<code>true</code> <code>false</code>	<i>values:</i> constant true constant false	

Figure 1.1: Booleans (B)

Denotational semantics In this approach developed by Scott and Strachey each program is interpreted as some mathematical object that belongs to *semantic domain*. Originally the domain was a space that included self-applicable functions and the meaning was obtained as a refinement (sequence of objects converging to result). One can talk in abstract terms about program behavior, e.g. program equivalence.

Operational semantics is defined in terms of *transition relations*. The evaluation of program consists in steps of *state* transformations. Thus evaluation of the terms can be thought as the work of *abstract machine*. This semantics is used throughout the notes.

1.1.4 Evaluation

Let us define the operational semantics for Boolean expressions only. The summary is on figure 1.1. Special subset of terms is separated as *values*⁴. The r.h.s. defines the *evaluation relation* on terms \longrightarrow . Some authors rather call it *reduction* relation. Definition is done in the same manner as with terms themselves, that is using inductive definition based on inference rules.

Definition 1.5. The *one-step* evaluation⁵ relation \longrightarrow is the smallest binary relation on terms satisfying the three rules in Figure 1.1. When the pair

⁴Ukrainian equivalent: *значення*

⁵Ukrainian equivalent: *один крок обчислення*

$(\mathfrak{t}, \mathfrak{t}')$ is in the evaluation relation, we say that “the evaluation *statement* (or *judgment*) $\mathfrak{t} \longrightarrow \mathfrak{t}'$ is *derivable*”.

The rules together determine how the terms are evaluated, i.e. evaluation *strategy*⁶. In this particular case evaluation is done from outermost to innermost conditionals, with inner conditionals being evaluated only for the sake of outermost guards. Rules E-IFTRUE, E-IFFALSE determine the evaluation for terms with fully evaluated guard. Rule E-IF extends the evaluation on other terms. Rules E-IFTRUE, E-IFFALSE are called *computational* rules, E-IF – *congruence* rule.

Proof that particular evaluation statement is derivable is usually done using (evaluation) derivation *tree*⁷. It is a structure similar to abstract syntax tree. Its nodes are labeled by inference rules and its root is labeled by evaluation statement derived by the tree.

If we abbreviate

$$\begin{aligned} u &= \text{if true then false else false} \\ t &= \text{if } u \text{ then } t_1 \text{ else } t_2 \\ s &= \text{if false then } t_1 \text{ else } t_2 \end{aligned}$$

Then we have derivation tree

$$\frac{\frac{\frac{}{u \longrightarrow \text{false}} \text{E-IFTRUE}}{t \longrightarrow s} \text{E-IF}}{\text{if } t \text{ then } t_3 \text{ else } t_4 \longrightarrow \text{if } s \text{ then } t_3 \text{ else } t_4} \text{E-IF}}$$

The derivation tree corresponding to evaluation statement is simply called its *derivation*⁸. The *evaluation* derivation⁹ degenerates to list because inference rules in example has at most one premise.

Similarly to induction on terms we can introduce proof technique called induction *on derivations*¹⁰.

Proposition 1.6 (Determinacy of one-step evaluation). *If $t \longrightarrow t'$ and $t \longrightarrow t''$, then $t' = t''$.*

⁶Ukrainian equivalent: *стратегія обчислення*

⁷Ukrainian equivalent: *дерево виведення*

⁸Ukrainian equivalent: *виведення*

⁹Ukrainian equivalent: *виведення (кроку) обчислення*

¹⁰Ukrainian equivalent: *індукція за виведеннями*

Proof. By induction on derivation of $t \longrightarrow t'$.¹¹ We have to consider cases according to rule applied last (at the root of derivation).

Case: Last rule is E-IFTRUE. Then $t = \text{if true then } t' \text{ else } t_1$. Now both E-IFFALSE, E-IF are not applicable to t and the only possible conclusion of E-IFTRUE is $t \longrightarrow t'$. Thus $t' = t''$.

Case: Last rule is E-IFFALSE. Similar to previous.

Case: Last rule in derivation of $t \longrightarrow t'$ is E-IF. Then t has the form $\text{if } t_1 \text{ then } t_2 \text{ else } t_3$, where $t_1 \longrightarrow t'_1$ for some t'_1 . This excludes $t_1 = \text{true}$, $t_1 = \text{false}$. The last rule in derivation of $t \longrightarrow t''$ then should be also E-IF, i.e. $t'' = \text{if } t''_1 \text{ then } t_2 \text{ else } t_3$, where $t_1 \longrightarrow t''_1$. By induction hypothesis $t'_1 = t''_1$, as a result $t' = t''$. \square

Definition 1.7. A term t is in *normal form*¹² if no evaluation rule applies to it, i.e. if there is no t' such that $t \longrightarrow t'$.

Trivially, following theorem holds.

Theorem 1.8. *Every value is in normal form.*

The converse to it also true.

Theorem 1.9. *If t is in normal form, then t is a value.*

Proof. Suppose that t is not a value. We're going to prove by induction on term structure that t is not a normal form. Either $t = \text{if true then } t' \text{ else } t''$. In which case by E-IFTRUE $t \longrightarrow t'$. Or $t = \text{if false then } t' \text{ else } t''$. Then by rule E-IFTRUE $t \longrightarrow t''$. Or, finally, $t = \text{if } t_1 \text{ then } t_2 \text{ else } t_3$, where t_1 is not a value. By inductive hypothesis there is term t'_1 such that $t_1 \longrightarrow t'_1$. Then by rule E-IF we conclude $t \longrightarrow \text{if } t'_1 \text{ then } t_2 \text{ else } t_3$. That is in all cases t is not a normal form. \square

Definition 1.10. The *multi-step evaluation* relation \longrightarrow^* is the reflexive, transitive closure of one-step evaluation.

Proposition 1.11 (Uniqueness of multi-step evaluation relation). *If $t \longrightarrow^* u$ and $t \longrightarrow^* u'$, where u and u' are both normal forms, then $u = u'$.*

Proof. It is a corollary of single step evaluation. \square

Theorem 1.12 (Termination of evaluation). *For every term t there is some normal form t' such that $t \longrightarrow^* t'$.*

¹¹This proof can be done by induction on term t and by induction on derivation of $t \longrightarrow t''$ as well

¹²Ukrainian equivalent: *нормальна форма*

New syntactic forms		New evaluation rules		$t \longrightarrow t'$
$t ::= \dots$	<i>terms:</i>	$\frac{t_1 \longrightarrow t'_1}{\text{succ } t_1 \longrightarrow \text{succ } t'_1}$		(E-SUCC)
0	constant zero	$\text{pred } 0 \longrightarrow 0$		(E-PREDZERO)
$\text{succ } t$	successor	$\text{pred } (\text{succ } nv_1) \longrightarrow nv_1$		(E-PREDSUCC)
$\text{pred } t$	predecessor	$\frac{t_1 \longrightarrow t'_1}{\text{pred } t_1 \longrightarrow \text{pred } t'_1}$		(E-PRED)
$\text{iszero } t$	zero test	$\text{iszero } 0 \longrightarrow \text{true}$		(E-ISZEROZERO)
$v ::= \dots$	<i>values:</i>	$\text{iszero } (\text{succ } nv_1) \longrightarrow \text{false}$		(E-ISZEROSUCC)
nv	numeric value	$\frac{t_1 \longrightarrow t'_1}{\text{iszero } t_1 \longrightarrow \text{iszero } t'_1}$		(E-ISZERO)
$nv ::= \dots$	<i>numeric values:</i>			
0	zero value			
$\text{succ } nv$	successor value			

Figure 1.2: Arithmetic expressions (NB)

Exercise 1.13. 1. Suppose we add a new rule

$$\text{if true then } t_2 \text{ else } t_3 \longrightarrow t_3 \quad (\text{E-FUNNY1})$$

to the ones in Figure 1.1. Which of the above theorems (3.5.4,7,8,11,12 uniqueness, termination) remain valid?

2. Suppose instead that we add this rule:

$$\frac{t_2 \longrightarrow t'_2}{\text{if } t_1 \text{ then } t_2 \text{ else } t_3 \longrightarrow \text{if } t_1 \text{ then } t'_2 \text{ else } t_3} \quad (\text{E-FUNNY2})$$

Which of the above theorems remain valid? Do any of the proofs need to change?

Let us now extend evaluation relation to numbers. See additional syntactic forms and evaluation rules in Figure 1.2.

There is a subset of terms called *numeric values*, which correspond to natural numbers. There are four computation rules and three congruence rules.

Definition 1.14. A term is *stuck*¹³ if it is in normal form but not a value.

¹³Ukrainian equivalent: *застряглий*

<i>Evaluation</i>	$\boxed{\tau \Downarrow v}$		
$v \Downarrow v$	(B-VALUE)		$\frac{\tau_1 \Downarrow 0}{\text{pred } \tau_1 \Downarrow 0}$ (B-PREDZERO)
$\frac{\tau_1 \Downarrow \text{true} \quad \tau_2 \Downarrow v_2}{\text{if } \tau_1 \text{ then } \tau_2 \text{ else } \tau_3 \Downarrow v_2}$	(B-IFTRUE)		$\frac{\tau_1 \Downarrow \text{succ } nv_1}{\text{pred } \tau_1 \Downarrow nv_1}$ (B-PREDSUCC)
$\frac{\tau_1 \Downarrow \text{false} \quad \tau_3 \Downarrow v_3}{\text{if } \tau_1 \text{ then } \tau_2 \text{ else } \tau_3 \Downarrow v_3}$	(B-IFFALSE)		$\frac{\tau_1 \Downarrow 0}{\text{iszero } \tau_1 \Downarrow \text{true}}$ (B-ISZEROZERO)
$\frac{\tau_1 \Downarrow nv_1}{\text{succ } \tau_1 \Downarrow \text{succ } nv_1}$	(B-SUCC)		$\frac{\tau_1 \Downarrow \text{succ } nv_1}{\text{iszero } \tau_1 \Downarrow \text{false}}$ (B-ISZEROSUCC)

Figure 1.3: Big-step semantics rules

Big-step semantics is an alternative way to define term evaluation. Unlike small-step semantics it is aimed to convey whole evaluation process in one step. The evaluation rules for arithmetics are provided in Figure 1.3

1.2 Typed arithmetic

Despite being quite simple the system of arithmetic expressions is enough to demonstrate most important aspects of type usage.

1.2.1 Typing relation

We separate rules for Booleans and for numbers. Details are in the figures 1.4 and 1.5 respectively.

Definition 1.15. The *typing relation*¹⁴ for arithmetic expressions is the smallest binary relation between terms and types satisfying all instances of the rules in Figures 1.4 and 1.5. A term τ is *typable*¹⁵ (or *well-typed*¹⁶) if there is some T such that $\tau : T$.

¹⁴Ukrainian equivalent: *відношення типізації*

¹⁵Ukrainian equivalent: *типізований*

¹⁶Ukrainian equivalent: *коректно типізований*

\mathbb{B} (typed)		<i>Extends</i> \mathbf{B}
<p><i>New syntactic forms</i></p> $\mathbf{T} ::= \mathbf{Bool}$ <p style="text-align: right;"><i>types:</i> <i>type of booleans</i></p>	<p><i>New typing rules</i></p> <div style="border: 1px solid black; padding: 2px; display: inline-block; margin-bottom: 5px;">$t : \mathbf{T}$</div> $\frac{}{\mathbf{true} : \mathbf{Bool}} \quad (\mathbf{T}\text{-TRUE})$ $\frac{}{\mathbf{false} : \mathbf{Bool}} \quad (\mathbf{T}\text{-FALSE})$ $\frac{t_1 : \mathbf{Bool} \quad t_2 : \mathbf{T} \quad t_3 : \mathbf{T}}{\mathbf{if } t_1 \text{ then } t_2 \text{ else } t_3 : \mathbf{T}} \quad (\mathbf{T}\text{-IF})$	

Figure 1.4: Typing rules for booleans (B)

\mathbb{N} (typed)		<i>Extends</i> \mathbf{NB}
<p><i>New syntactic forms</i></p> $\mathbf{T} ::= \dots$ <p style="text-align: right;"><i>types:</i> <i>type of natural numbers</i></p>	<p><i>New typing rules</i></p> <div style="border: 1px solid black; padding: 2px; display: inline-block; margin-bottom: 5px;">$t : \mathbf{T}$</div> $\frac{}{0 : \mathbf{Nat}} \quad (\mathbf{T}\text{-ZERO})$	$\frac{t_1 : \mathbf{Nat}}{\mathbf{succ } t_1 : \mathbf{Nat}} \quad (\mathbf{T}\text{-SUCC})$ $\frac{t_1 : \mathbf{Nat}}{\mathbf{pred } t_1 : \mathbf{Nat}} \quad (\mathbf{T}\text{-PRED})$ $\frac{t_1 : \mathbf{Nat}}{\mathbf{iszero } t_1 : \mathbf{Bool}} \quad (\mathbf{T}\text{-ISZERO})$

Figure 1.5: Typing rules for numbers (NB)

Similar to evaluation relation typing of some term is proved using *typing* derivation¹⁷.

Here is the typing derivation for the term `if iszero 0 then 0 else pred 0`:

$$\frac{\frac{}{0 : \mathbf{Nat}} \quad \mathbf{T}\text{-ZERO} \quad \frac{}{\mathbf{iszero } 0 : \mathbf{Bool}} \quad \mathbf{T}\text{-ISZERO}}{0 : \mathbf{Nat}} \quad \mathbf{T}\text{-ZERO} \quad \frac{}{0 : \mathbf{Nat}} \quad \mathbf{T}\text{-ZERO} \quad \frac{}{\mathbf{pred } 0 : \mathbf{Nat}} \quad \mathbf{T}\text{-PRED}}{\mathbf{if iszero } 0 \text{ then } 0 \text{ else } \mathbf{pred } 0 : \mathbf{Nat}} \quad \mathbf{T}\text{-IF}$$

Simplest property of the type system is presented by the following lemma.

Lemma 1.16 (Inversion Lemma). *1. If `true` : R, then R = Bool.*

¹⁷Ukrainian equivalent: *виведення типу*

2. If $false : R$, then $R = Bool$.
3. If $if\ t_1\ then\ t_2\ else\ t_3 : R$, then $t_1 : Bool$, $t_2 : R$ and $t_3 : R$.
4. If $0 : R$, then $R = Nat$.
5. If $succ\ t_1 : R$, then $R = Nat$ and $t_1 : Nat$.
6. If $pred\ t_1 : R$, then $R = Nat$ and $t_1 : Nat$.
7. If $iszero\ t_1 : R$, then $R = Bool$ and $t_1 : Nat$.

Proof. All points are quite similar. Let us show one of them. If term of the form $if\ t_1\ then\ t_2\ else\ t_3$ is typable, then the last rule of typing derivation should be T-IF. Therefore $t_1 : Bool$, $t_2 : R$, $t_3 : R$. \square

Theorem 1.17 (Uniqueness of Types). *Each term t has at most one type. That is, if t is typable, then its type is unique.*

Proof. By induction on typing derivation of $t : T$. The only nontrivial case is T-IF. The term $t = if\ t_1\ then\ t_2\ else\ t_3$ is only typable when the last rule in typing derivation is T-IF (inversion lemma). By i.h. type of t_1 and t_2 is unique. Therefore type of t is also unique. \square

1.2.2 Type safety

*Type safety*¹⁸ should allow us to avoid meaningless (wrong) terms. This can be expressed in two properties of type system.

*Progress*¹⁹: A well-typed term is not stuck (either it is a value or it can take a step according to the evaluation rules).

*Preservation*²⁰: If a well-typed term takes a step of evaluation, then the resulting term is also well typed.

In order to discuss and prove these properties the following lemma is required.

Lemma 1.18 (Canonical Forms). *1. If v is a value of type $Bool$, then v is either $true$ or $false$.*

2. If v is a value of type Nat , then v is a numeric value according to the grammar.

¹⁸Ukrainian equivalent: *типова безпека*

¹⁹Ukrainian equivalent: *просування*

²⁰Ukrainian equivalent: *збереження*

Proof. By inversion lemma numeric values cannot have type `Bool` and vice versa. \square

Theorem 1.19 (Progress). *Suppose t is a well-typed term (that is, $t : T$ for some T). Then either t is a value or else there is some t' with $t \longrightarrow t'$.*

Proof. By induction on typing derivation of $\mathfrak{t} : T$. We consider only some cases.

Case T-TRUE. \mathfrak{t} is a value in this case.

Case T-PRED. The term has form `pred \mathfrak{t}_1` where $\mathfrak{t}_1 : \text{Nat}$. If \mathfrak{t}_1 is a value then by canonical form lemma $\mathfrak{t}_1 = \text{nv}_1$. Therefore E-PREDZERO or E-PREDSUCC is applicable to \mathfrak{t} . Otherwise by i.h. $\mathfrak{t}_1 \longrightarrow \mathfrak{t}'_1$ then by E-PRED $\mathfrak{t} \longrightarrow \text{pred } \mathfrak{t}'_1$. \square

Theorem 1.20 (Preservation (*subject reduction*²¹)). *If $t : T$ and $t \longrightarrow t'$, then $t' : T$.*

1.3 Untyped λ -calculus

Historical background Historically lambda-calculus notation was inspired by the notation used by Russell and Whitehead in Principia Mathematica[6]. Russell used expressions with circumflex like $\hat{x}(\phi x)$. Roughly speaking, this (class abstraction) represents the set denoted $\{x \mid \phi x\}$ in traditional set-builder notation[7]. Also at some point to some extent it was used to express functions $\phi\hat{x}$. For metaphysical reasons Russell abandoned treatment of functions as entities and further discarded the notation in favour of substitutional theory.

Because of typesetting problems with circumflex over letter, Church first used circumflex in front of variable and then for similarity with Λ started using letter λ for abstraction.

1.3.1 Lambda terms

Understanding of free and bound variables. See: wiki.

Let us introduce function $FV(\mathfrak{t})$ that gives free variable in term \mathfrak{t} . By induction on term structure:

$$\begin{aligned} FV(\mathfrak{x}) &= \{\mathfrak{x}\} \\ FV(\lambda \mathfrak{x}. \mathfrak{t}) &= FV(\mathfrak{t}) \setminus \{\mathfrak{x}\} \\ FV(\mathfrak{t}_1 \mathfrak{t}_2) &= FV(\mathfrak{t}_1) \cup FV(\mathfrak{t}_2) \end{aligned}$$

²¹Ukrainian equivalent: *скорочення/редукція підмета*

λ (α -equivalence part)

<p><i>Evaluation</i></p> <div style="border: 1px solid black; padding: 5px; margin: 10px auto; width: 80px; text-align: center;"> $t \equiv t$ </div> <div style="border: 1px solid black; padding: 5px; margin: 10px auto; width: 150px; text-align: center;"> $t \equiv t' \quad t' \equiv t''$ $t \equiv t''$ </div>	<div style="border: 1px solid black; padding: 2px; margin: 0 auto; width: 60px; text-align: center;"> $t \overset{\alpha}{\equiv} t'$ </div>	<p>(A-REF)</p>	<div style="border: 1px solid black; padding: 5px; margin: 10px auto; width: 150px; text-align: center;"> $t_1 \equiv t'_1 \quad t_2 \equiv t'_2$ $t_1 t_2 \equiv t'_1 t'_2$ </div>	<p>(A-APP)</p>
<div style="border: 1px solid black; padding: 5px; margin: 10px auto; width: 150px; text-align: center;"> $t \equiv t' \quad t' \equiv t''$ $t \equiv t''$ </div>	<p>(A-TRANS)</p>	<div style="border: 1px solid black; padding: 5px; margin: 10px auto; width: 150px; text-align: center;"> $y \notin FV(t)$ $\lambda x. t \equiv \lambda y. [x \mapsto y]t$ </div>	<p>(A-ABS)</p>	

Figure 1.6: α -equivalence

Consider two expressions

$$\begin{aligned} \text{id}(x) &= x \\ \text{id}(y) &= y \end{aligned}$$

They can be considered both as calculations and as definition of id . It is obvious that defined function is the same and from first equation we derive $\text{id} = \lambda x. x$ while from the second $\text{id} = \lambda y. y$. This highlights the property that the names of the bound variables is not important and that bound variables can be changed without change of meaning. This is formalized by the following transformation. We define the relation \equiv_α called α -equivalence. Its inference rules are presented in figure 1.6.

Now we can define substitution rules by induction on term structure:

$$\begin{aligned} [x \mapsto t]x &= t \\ [x \mapsto t]y &= y \\ [x \mapsto t](t_1 t_2) &= [x \mapsto t]t_1 [x \mapsto t]t_2 \\ [x \mapsto t]\lambda x. t_1 &= \lambda x. t_1 \end{aligned}$$

Lambda abstraction is trickier. The abstraction by substituted variable shield the term. If the variable is different we have to make sure that substitution avoids *name capturing*, i.e. free variables of term t_1 should not become bound after substitution.

$$\frac{y \notin FV(t)}{[x \mapsto t]\lambda y. t_1 = \lambda y. [x \mapsto t]t_1}$$

Last rule is not restrictive because we can guarantee the premise for term that is α -equivalent to given term s thus extending substitution to all terms.

Notice that α -equivalence and substitution are defined by *mutual* induction (recursion).

Exercise 1.21. Prove that relation \equiv^α is an equivalence relation.

Exercise 1.22. Prove that following rule

$$[x \mapsto t]\lambda x. t_1 = \lambda x. t_1$$

can be eliminated from the substitution definition.

Exercise 1.23. Prove that function $FV(t)$ is correctly defined w.r.t. relation \equiv^α .

Further, as usual in λ -calculus, equality of terms is understood up to α -equivalence.

1.3.2 Beta reduction

The main element of λ -calculus semantics is expression of the kind $(\lambda x. t) t_2$, which is called *reducible expression*²² or *redex*²³ for short. The intuition behind the application of function to the given argument is to replace the abstracted variable in the expression of the function with t_2 . This is captured by the rule:

$$(\lambda x. t) t_2 \longrightarrow [x \mapsto t_2]t$$

Here $[x \mapsto t_2]$ is a substitution operator that will be formalized later in the section.

1.3.2.1 Evaluation strategies

Evaluation strategy determines in which order expressions are to be evaluated. Since redexes are in the core of lambda-calculus evaluation we should choose which one to evaluate first. Additional rules guide us through the expression. For lambda-calculus enriched with additional elements they also regulate is it redex reduction or specific rules have the preference in evaluation.

- full beta-reduction
- normal order

²²Ukrainian equivalent: *скорочуваний вираз, редукований вираз*

²³Ukrainian equivalent: *редекс*

- call by name / call by need
- call by value

Mostly last two strategies are implemented in programming languages. Imperative languages use call-by-value. Functional programming language are divided. Some like OCaml or prefer call-by-value, it is called *eager*²⁴, *strict*²⁵ or *greedy*²⁶. Haskell uses call by need strategy, which is also called *lazy* evaluation²⁷.

1.3.2.2 Other rules

For some systems additional η -conversion rule is also considered. It written as follows

$$\lambda x.(f x) \rightarrow f$$

To be valid there is a constraint that does not contain x free. It represents to extensionality. Depending on the direction in which it is applied we could have either η -reduction (direct) or η -expansion (reverse). Addition of this rule allows to perform certain optimization during compilation. But it also means adding complexity in compiler.

1.3.3 Programming in pure λ -calculus

λ -calculus is proven to be Turing-complete. Arithmetic, partial recursive functions could be implemented in it. To be accessible from λ -calculus notions from different mathematical systems require specific representation as λ -terms.

There are “canonical” representation of well-known mathematical constructions in pure λ -calculus. There are booleans, pairs, (Church) numbers. Basic arithmetic operation of addition, multiplication and exponentiation are implemented straightforward. More complex constructions require recursion. It is implemented in λ -calculus using fixed-point combinator. Good account of these implementations is presented in [2].

Exercise 1.24. Implement in pure lambda-calculus iteration without fixed-point combinator.

Exercise 1.25. Implement in enriched lambda-calculus pairs based on embedded booleans.

²⁴Ukrainian equivalent: *жадібне*

²⁵Ukrainian equivalent: *строге*

²⁶Ukrainian equivalent: *жадібне*

²⁷Ukrainian equivalent: *ліниві* обчислення

\rightarrow (untyped)							
<p><i>Syntax</i></p> <p>$\mathbf{t} ::=$</p> <div style="margin-left: 20px;"> \mathbf{x} $\lambda \mathbf{x}. \mathbf{t}$ $\mathbf{t} \ \mathbf{t}$ </div> <p>$\mathbf{v} ::=$</p> <div style="margin-left: 20px;"> $\lambda \mathbf{x}. \mathbf{t}$ </div>	<p><i>terms:</i></p> <p><i>variable</i></p> <p><i>abstraction</i></p> <p><i>application</i></p> <p><i>values:</i></p> <p><i>abstraction value</i></p>						
<div style="border: 1px solid black; display: inline-block; padding: 2px 5px;">$\mathbf{t} \longrightarrow \mathbf{t}'$</div>							
<p><i>Evaluation</i></p> <div style="margin-left: 100px;"> <table style="border-collapse: collapse;"> <tr> <td style="border: 1px solid black; padding: 5px; text-align: center;"> $\frac{\mathbf{t}_1 \longrightarrow \mathbf{t}'_1}{\mathbf{t}_1 \ \mathbf{t}_2 \longrightarrow \mathbf{t}'_1 \ \mathbf{t}_2}$ </td> <td style="padding-left: 20px; vertical-align: middle;">(E-APP1)</td> </tr> <tr> <td style="border: 1px solid black; padding: 5px; text-align: center;"> $\frac{\mathbf{t}_2 \longrightarrow \mathbf{t}'_2}{\mathbf{v}_1 \ \mathbf{t}_2 \longrightarrow \mathbf{v}_1 \ \mathbf{t}'_2}$ </td> <td style="padding-left: 20px; vertical-align: middle;">(E-APP2)</td> </tr> <tr> <td style="border: 1px solid black; padding: 5px; text-align: center;"> $(\lambda \mathbf{x}. \mathbf{t}_{12}) \ \mathbf{v}_2 \longrightarrow [\mathbf{x} \mapsto \mathbf{v}_2] \mathbf{t}_{12}$ </td> <td style="padding-left: 20px; vertical-align: middle;">(E-APPABS)</td> </tr> </table> </div>		$\frac{\mathbf{t}_1 \longrightarrow \mathbf{t}'_1}{\mathbf{t}_1 \ \mathbf{t}_2 \longrightarrow \mathbf{t}'_1 \ \mathbf{t}_2}$	(E-APP1)	$\frac{\mathbf{t}_2 \longrightarrow \mathbf{t}'_2}{\mathbf{v}_1 \ \mathbf{t}_2 \longrightarrow \mathbf{v}_1 \ \mathbf{t}'_2}$	(E-APP2)	$(\lambda \mathbf{x}. \mathbf{t}_{12}) \ \mathbf{v}_2 \longrightarrow [\mathbf{x} \mapsto \mathbf{v}_2] \mathbf{t}_{12}$	(E-APPABS)
$\frac{\mathbf{t}_1 \longrightarrow \mathbf{t}'_1}{\mathbf{t}_1 \ \mathbf{t}_2 \longrightarrow \mathbf{t}'_1 \ \mathbf{t}_2}$	(E-APP1)						
$\frac{\mathbf{t}_2 \longrightarrow \mathbf{t}'_2}{\mathbf{v}_1 \ \mathbf{t}_2 \longrightarrow \mathbf{v}_1 \ \mathbf{t}'_2}$	(E-APP2)						
$(\lambda \mathbf{x}. \mathbf{t}_{12}) \ \mathbf{v}_2 \longrightarrow [\mathbf{x} \mapsto \mathbf{v}_2] \mathbf{t}_{12}$	(E-APPABS)						

Figure 1.7: Untyped lambda-calculus (λ)

1.3.4 Formal definitions

The summary of the rules is presented in the figure 1.7.

It is easy to see there what strategy is implied here. Notice that λ -abstractions are values. The term representing the function is evaluated first until it is a value. Then its argument is evaluated. Finally, the redex itself is reduced.

1.4 Simply-typed λ -calculus

Next system is known as *simply-typed* lambda calculus²⁸. The summary of the rules is presented in the figure 1.8. Its set of types is inductively defined using several base types and *function* type²⁹ constructor. Type annotations appear for variables in λ -abstractions. Context Γ is introduced. It is needed to assign types to free variables so that open terms could be typed. Typing rules express that function type and its parameter type must agree. Also it specify that function type is assigned to the λ -abstraction.

1.4.1 Properties of λ_{\rightarrow}

Lemma 1.26 (Canonical form). *1. If v is a value of type $Bool$, then v is either *true* or *false*.*

2. If v is a value of type $T_1 \rightarrow T_2$, then $v = \lambda x : T_1. t_2$.

²⁸Ukrainian equivalent: лямбда-числення з простою типізацією

²⁹Ukrainian equivalent: “стрілочний” тип

\rightarrow (typed)Based on λ

Syntax		Evaluation		Typing	
$t ::=$	x $\lambda x : T . t$ $t t$	<i>terms:</i> <i>variable</i> <i>abstraction</i> <i>application</i>	$t \longrightarrow t'$		
$v ::=$	$\lambda x : T . t$	<i>values:</i> <i>abstraction value</i>	$\frac{t_1 \longrightarrow t'_1}{t_1 t_2 \longrightarrow t'_1 t_2}$	(E-APP1)	
$T ::=$	$T \rightarrow T$	<i>types:</i> <i>type of functions</i>	$\frac{t_2 \longrightarrow t'_2}{v_1 t_2 \longrightarrow v_1 t'_2}$	(E-APP2)	
$\Gamma ::=$	\emptyset $\Gamma, x : T$	<i>contexts:</i> <i>empty context</i> <i>term variable binding</i>	$(\lambda x : T_{11} . t_{12}) v_2 \longrightarrow [x \mapsto v_2] t_{12}$	(E-APPABS)	
					$\Gamma \vdash t : T$
			$\frac{x : T \in \Gamma}{\Gamma \vdash x : T}$	(T-VAR)	
			$\frac{\Gamma, x : T_1 \vdash t_2 : T_2}{\Gamma \vdash \lambda x : T_1 . t_2 : T_1 \rightarrow T_2}$	(T-ABS)	
			$\frac{\Gamma \vdash t_1 : T_{11} \rightarrow T_{12} \quad \Gamma \vdash t_2 : T_{11}}{\Gamma \vdash t_1 t_2 : T_{12}}$	(T-APP)	

Figure 1.8: Pure simply typed lambda-calculus (λ_{\rightarrow})

Proposition 1.27 (Progress). *Suppose t is a closed, well-typed term (that is, $\vdash t : T$ for some T). Then either t is a value or else there is some t' with $t \longrightarrow t'$.*

Proof. By induction on typing derivation of $\vdash t : T$. The proof has two parts: for pure λ -calculus and for booleans. Proof for booleans remain the same. Here we consider proof for pure λ_{\rightarrow} . It is obvious that t is not a variable, so case T-VAR is vacuously true. If the last rule is T-ABS then t is immediately a value. Thus only case of application should be considered.

Case T-APP: t has a form $t_2 t_1$ with $\vdash t_2 : T_1 \rightarrow T$ and $\vdash t_1 : T_1$.

Then by inductive hypothesis for $\vdash t_2 : T_1 \rightarrow T$ either t_2 is a value or evaluates in one step to t'_2 . If $t_2 \longrightarrow t'_2$ then by E-APP1 $t \longrightarrow t'_2 t_1$. Otherwise by canonical form lemma t_2 should have a form $\lambda x : T_1 . t_3$, i.e. t is a redex. By inductive hypothesis for $\vdash t_1 : T_1$ either t_1 is a value or $t_1 \longrightarrow t'_1$. In former case E-APPABS applies to t , and in the latter E-APP2 does. \square

Lemma 1.28 (Preservation of types under substitution). *If $\Gamma, x:S \vdash t : T$ and $\Gamma \vdash s : S$ then $\Gamma \vdash [x \mapsto s]t : T$.*

Proof. By induction on typing derivation of $\Gamma \vdash t : T$.

Case T-VAR.

Subcase 1. $t = y, \Gamma \vdash y : T, y:T \in \Gamma$. Then $[x \mapsto s]y = y$, and respectively $\Gamma \vdash [x \mapsto s]y : T$.

Subcase 2. $t = x, T = S$. Then $[x \mapsto s]x = s$, and $\Gamma \vdash [x \mapsto s]x : T$.

Case T-ABS. $t = \lambda y:T_1. t_2, T = T_1 \rightarrow T_2, \Gamma, x:S, y:T_1 \vdash t_2 : T_2$. Due to α -equivalence we can assume that $y \notin \{x\} \cup FV(s)$. Therefore $[x \mapsto s]t = \lambda y:T_1. [x \mapsto s]t_2$. Reordering last two elements of typing context for t_2 we obtain $\Gamma, y:T_1, x:S \vdash t_2 : T_2$. Also we can weaken the typing of $s: \Gamma, y:T_1 \vdash s : S$. By inductive hypothesis $y:T_1, \Gamma \vdash [x \mapsto s]t_2 : T_2$. Then by rule T-ABS $\Gamma \vdash \lambda y:T_1. [x \mapsto s]t_2 : T_1 \rightarrow T_2$.

Case T-APP. $t = t_1 t_2, \Gamma, x:S \vdash t_1:T_2 \rightarrow T, \Gamma, x:S \vdash t_2:T_2$.

Then by inductive hypothesis $\Gamma \vdash [x \mapsto s]t_2 : T_2$ and $\Gamma \vdash [x \mapsto s]t_1 : T_2 \rightarrow T$. By T-APP, $\Gamma \vdash [x \mapsto s](t_1 t_2) : T$ \square

Proposition 1.29 (Preservation). *If $\Gamma \vdash t : T$ and $t \longrightarrow t'$, then $\Gamma \vdash t' : T$.*

Proof. By induction on typing derivation of $\Gamma \vdash t : T$.

Cases T-VAR, T-ABS are vacuously true because their conclusion term is a normal form.

Case T-APP: t has a form $t_2 t_1$ with $\Gamma \vdash t_2 : T_1 \rightarrow T$ and $\Gamma \vdash t_1 : T_1$.

We consider subcases according to last rule of evaluation derivation

Subcase E-APP1: $t_2 \longrightarrow t'_2$

Then t' is $t'_2 t_1$. By inductive hypothesis $\Gamma \vdash t'_2 : T_1 \rightarrow T$ and $t' : T$ by rule T-APP.

Subcase E-APP2: $t_2 = v_2$ is a value and $t_1 \longrightarrow t'_1$

Then t' is $v_2 t'_1$. By inductive hypothesis $\Gamma \vdash t'_1 : T_1$. Therefore by rule T-APP $\Gamma \vdash t' : T_2$.

Subcase E-APPABS: $(\lambda x:T_1. t_3) v_1$ as t

This is a redex that evaluates to $t' = [x \mapsto v_1]t_3$. Since $\Gamma, x:T_1 \vdash t_3 : T$ and $\Gamma \vdash v_1 : T_1$ then by substitution lemma $\Gamma \vdash [x \mapsto v_1]t_3 : T$. \square

1.4.2 Type reconstruction problem

It has already been shown that it is typical to start with untyped system or language. It is easier to implement and it is more efficient in program evalu-

ation. Then in order to provide safety the language is augmented with types. In this setting question arise if it is possible to take program from original system (without types) and construct from it typed version with identical runtime behavior. This problem is called *type reconstruction*³⁰ problem. Besides being peculiar theoretical problem it has straightforward application: construction tool can assist programmer in identifying types of variables. This is especially viable when such types are quite involved, which is common for modern libraries, typical example are iterator types in C++. In this more local sense the problem of automatic type detection of expression based on its context is referred to as *type inference*³¹. This allows us to use advantages in safety of typed system and to avoid type annotations that seem redundant or get rid of them completely.

Type reconstruction can be considered as the inverse process to *type erasure*³². Like type reconstruction type erasure connects typed and untyped language. It is the process of removing information about types from the program before executing it. The most important property of type erasure is that it commutes with evaluation. But it also establishes close relationship between untyped and typed languages.

In the simplest case of typed arithmetic there are no type annotations at all therefore type erasure is identity. Semantics for typed arithmetic is implicitly preserved because it is based on the same evaluation rules. Therefore typed arithmetic term are evaluated as expected. Typed reconstruction in this case is type-checking that is finding out if the term is typable.

Let us consider simply-typed lambda calculus λ_{\rightarrow} . Erasure of a simply-typed term boils down to erasing type annotations. This can be expressed inductively as follows

$$\begin{aligned} \text{erase}(x) &= x \\ \text{erase}(\lambda x:T. t) &= \lambda x. \text{erase}(t) \\ \text{erase}(t_1 t_2) &= \text{erase}(t_1) \text{erase}(t_2) \end{aligned}$$

Since evaluation rules for simply-typed lambda terms and untyped lambda-terms are basically the same (up to type annotations), it is to see following proposition.

Proposition 1.30 (Erasure commutes with evaluation). *Following holds*

- If $t \rightarrow t'$ then $\text{erase}(t) \rightarrow \text{erase}(t')$

³⁰Ukrainian equivalent: *відновлення типів*

³¹Ukrainian equivalent: *виведення типів*

³²Ukrainian equivalent: *стирання типів*

- If $\text{erase}(t) \rightarrow m'$ then there is λ_{\rightarrow} -term t' such that $t \rightarrow t'$ and $m' = \text{erase}(t')$.

Now we can refine the notion of type reconstruction problem

Definition 1.31. Term m of untyped lambda-calculus is *typable* in λ_{\rightarrow} if there are t, T, Γ such that $\Gamma \vdash t : T$, $\text{erase}(t) = m$.

This concepts play important role in following chapter when we consider how polymorphic constructions can be added to the system.

Chapter 2

Polymorphism

2.1 Introduction

In conventional languages functions are monomorphic. Still there are examples of symbols that has multiple meaning. Christopher Strachey called this symbols *polymorphic*¹:

We expect compiling system to interpret correctly ambiguous symbols (such as +) which mean different things according to the types of their operands. We call ambiguous operators of this sort **polymorphic** as they have several forms depending on their arguments.[8]

Strachey divided such symbols into two categories. Some coincide by chance, other work uniformly. First situation was called *ad hoc* polymorphism². In this case polymorphic symbols work on several different types but have unrelated behaviour. Such polymorphic symbols have meaning only for limited number of cases. There is no single systematic way of determining the type of the result from the type of the arguments[8]. The second is called *parametric* polymorphism³. Parametric is regular and should work for any types in uniform way.

The intrinsic problem of polymorphism choosing correct form of the symbol from the context. The choice is complicated by the fact that several forms can be accepted for some contexts. Therefore

Luca Cardelli and Peter Wegner further refined the classification of polymorphism [1]. This classification is represented in figure 2.1. One can see

¹Ukrainian equivalent: *поліморфні, “багатоформні”*

²Ukrainian equivalent: *ad hoc* поліморфізм

³Ukrainian equivalent: *параметричний* поліморфізм

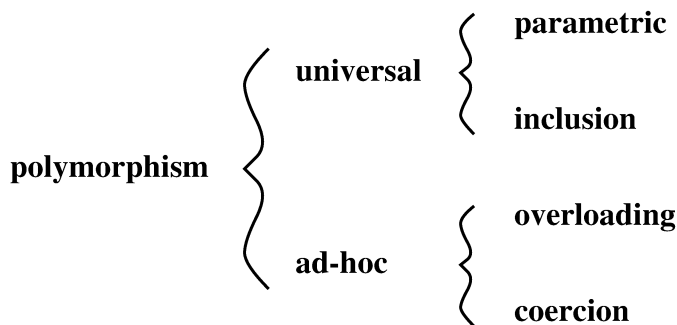


Figure 2.1: Polymorphism classification[1]

that parametric polymorphism when different types are treated uniformly is generalized to *universal* polymorphism⁴, while ad hoc polymorphism is further subdivided into overloading and *coercion*⁵. Let us briefly consider each of the refined categories.

2.1.1 Ad hoc polymorphism

*Overloading*⁶ is the ability to define several objects (functions) with the same name. It (degeneration) can be (and sometimes is) split by source code preprocessing stage. When used properly it is helpful tool that facilitates writing of the better code. Coercion is instead a semantic operation which is needed to convert an argument to the type expected by a function, in a situation which would otherwise result in a type error[1]. Strachey called such functions transfer functions. This could work either statically when application to arguments is generated at compile time, or decision about coercion specifics is made dynamically dictated by run-time evaluation of the arguments. The overloading and coercion can be combined in several ways. Typically overloading covers the limited list of most important cases, while intermediate cases are reduced to this list via coercion.

Monomorphic languages constrain object to one particular behavior which is too restrictive for expressive power. To lift this limitations polymorphism should be used. Therefore polymorphism is present in older programming languages. But they view polymorphic constructions as exceptions. They usually allow it in the following cases:

1. overloading

⁴Ukrainian equivalent: *універсальний поліморфізм*

⁵Ukrainian equivalent: *неявне приведення типів*

⁶Ukrainian equivalent: *довизначення, перевантаженя*

2. coercion
3. subtyping (sub/super ranges)
4. value sharing, like `nil`

The last two can be viewed as a special cases of subtyping and parametric polymorphism respectively.

Ad hoc polymorphism is also called *apparent* polymorphism⁷ because it disappears in the close range. Indeed for overloading several (possibly unrelated) monomorphic functions are named the same, and for coercion type of function result is fixed and independent of the argument type.

2.1.2 Subtype polymorphism

Languages with *subtype* polymorphism⁸ or subtyping provide the notion of *subtype*⁹. Type S is a subtype of T , denoted as $S <: T$, if terms of type S can be safely treated (be used in place) as terms of type T . This way any term belongs to multiple types at once and language, in fact, supports polymorphism. Many popular programming languages provide this kind of polymorphism. Often subtyping relation is based on the inheritance relation. For this reason subtyping is sometimes called *inheritance* polymorphism¹⁰. Due to requirement in some languages to include in subtype all methods of the supertype it is also called *inclusion* polymorphism¹¹. Most type systems equipped with subtyping have this relation to be a preorder.

One of the major differences of subtyping from parametric polymorphism is that code executed might depend on the particular type. This could add flexibility. At the same time, subtyping has intrinsic limitations. This is known as binary method problem[9].

Assume we want to manipulate in the program with the monoid algebraic structures. For this reason we introduce corresponding abstract class

```
interface Monoid {
    Monoid idElt();
    Monoid binop(Monoid y);
}
```

Now for particular examples of monoids we provide concrete classes. For instance:

⁷Ukrainian equivalent: *удаваний* поліморфізм

⁸Ukrainian equivalent: *підтиповий* поліморфізм

⁹Ukrainian equivalent: *підтип*

¹⁰Ukrainian equivalent: *спадковий* поліморфізм

¹¹Ukrainian equivalent: *включний* поліморфізм

```

class IntAddMonoid implements Monoid {
    public IntAddMonoid(int x) { n = x; }
    public IntAddMonoid binop(IntAddMonoid other)
        { return new IntAddMonoid(n + other.n); }

    public IntAddMonoid idElt()
        { return new IntAddMonoid(0); }

    public int n;
};

class DoubleAddMonoid implements Monoid {
    public DoubleAddMonoid(double x) { n = x; }

    public DoubleAddMonoid binop(DoubleAddMonoid other)
        { return new DoubleAddMonoid(n + other.n); }

    public DoubleAddMonoid idElt()
        { return new DoubleAddMonoid(0.0); }

    public double n;
};

```

Method binop signature is chosen according to method intended usage. Unfortunately this attempt fails with the following error message:

```

error: IntAddMonoid is not abstract and does not override abstract method
binop(Monoid) in Monoid

```

On the other hand if we try a workaround using type down cast:

```

class IntAddMonoid implements Monoid {
    public IntAddMonoid(int x) { n = x; }
    public IntAddMonoid binop(Monoid other)
        { return new IntAddMonoid(n + ((IntAddMonoid)other).n); }

    public IntAddMonoid id_elt()
        { return new IntAddMonoid(0); }
    int n;
};

```

Then it results in the following runtime error as soon as we pass non-IntAddMonoid instance as an argument to binop.

This situation is the consequence of subtyping rules. We want to implement the method of the abstract class therefore the type of the method in the subclass has to be compatible with type of this method in abstract

class. We come to the question: when $T' \rightarrow S' <: T \rightarrow S$? The requirements of subtyping imply that $S' <: S$ and $T <: T'$ in this case. From the standpoint of order theory function type constructor “ $_ \rightarrow _$ ” is monotone (order-preserving) by the second argument and antitone (order-reversing) by the first argument. Therefore argument type is called *contravariant*¹², while result type *covariant*¹³. Both of these terms stem from the category theory, where (among other examples) ordered sets are considered as categories and maps between them as functors.

In general any operator on types (*type constructor*) can be considered. For instance, type constructor for arrays:

$$\text{Array: } \mathcal{T} \rightarrow \mathcal{T}.$$

2.1.3 Introduction to parametric polymorphism

Parametric polymorphism is considered the purest form of polymorphism. As was mentioned earlier in this case function should work uniformly for any given types. This literally means that the same code is used independently of the argument types. Canonical example of this is the identity function¹⁴:

```
- fun id x = x;
val id = fn : 'a -> 'a
- id 1;
val it = 1 : int
- id 1.0;
val it = 1.0 : real
```

Because of the parametric polymorphism universality special measures and approaches should be taken. In order to pass arguments of various types uniformity of types representation and use of certain calling convention might be required. Also to accomplish anything nontrivial “hands-off” approach is utilized[10]. Let us demonstrate it with the example of function `map`:

```
- fun map f [] = []
-   | map f (x::ls) = (f x)::map f ls;
val map = fn : ('a -> 'b) -> 'a list -> 'b list

- fun dub x = x + x;
val dub = fn : int -> int
```

¹²Ukrainian equivalent: *контраваріантність*

¹³Ukrainian equivalent: *коваріантність*

¹⁴implemented in ML

```
- map dub [1,2,3];
val it = [2, 4, 6] : int list
```

You can see that this function makes no assumptions about type parameters. It relies on list structure using pattern matching. When it needs to touch elements of the list it uses extra parameter (function `f`) to accomplish this.

Some typesystems also support setting bounds to type parameters used in polymorphic construction. Usually these bounds are provided in the form of subtyping constraints. This case of parametric polymorphism is known as *bounded* polymorphism¹⁵. It is supported in Haskell with type classes, in C# and Java with generics, in Scala with either virtual types or generics. This approach is further developed in *generic programming*¹⁶, where data types are grouped using concepts. Concepts are defined in terms of requirements to the types that adhere to the concepts. There are two types of requirements syntactic (like having a certain signature), semantic (datatype satisfies certain axioms).

Finally, there is a question how freely polymorphic constructions can be used. This means whether operations with them are restricted or polymorphic entities can be used freely in all kinds of operations. Such entities got the name of *first-class citizens*¹⁷[8]. Consider the following small example of imaginary conversation with system similar to ML

```
- fun f g = (g 1, g 1.0)
- f id
```

First line defines function that takes polymorphic function as an argument which is further applied to arguments of two different types. Then function `f` is applied for identity. For this construction to be valid it should be possible to pass polymorphic functions as parameters. This is possible when system considers such functions as first-class citizens. Therefore polymorphism that allows such things is often called *first-class* polymorphism¹⁸. In reality given example fails with following output:

```
stdIn:1.12-1.24 Error: operator and operand do not agree
[overload - bad instantiation]
operator domain: 'Z[INT]
operand:         real
in expression:
  g 1.0
```

¹⁵Ukrainian equivalent: поліморфізм з обмеженнями

¹⁶Ukrainian equivalent: параметризоване програмування

¹⁷Ukrainian equivalent: об'єкти першого класу

¹⁸Ukrainian equivalent: поліморфізм першого класу

Reasons for this would be discussed in the following subsections.

2.2 Let-polymorphism

In this subsection we consider one variant of parametric polymorphism in detail. We would target the both sides of the question, i.e. formal specification of the system in the terms already used in the notes as well as implementation as practical algorithm.

We start with simply-typed lambda calculus λ_{\rightarrow} and through type reconstruction problem gradually accommodate it to “polymorphic environment”. The major departure from system λ_{\rightarrow} is extension of basic types with type variables (parameters): Types:

$$T ::= \text{BasicTypes} \quad | \quad \boxed{X_i} \quad | \quad T \rightarrow T$$

Next step is define manipulations with type variables:

Definition 2.1. *Type substitutions* σ is a finite mapping from type variables to types.

For example, $\sigma = [X \mapsto U, Y \mapsto T]$. Application to terms and contexts respectively:

$$\begin{aligned} \sigma(X) &= \begin{cases} T & \text{if } (X \mapsto T) \in \sigma \\ X & \text{else } (X \notin \text{dom}\sigma) \end{cases} \\ \sigma(\text{Bool}) &= \text{Bool} \\ \sigma(T_1 \rightarrow T_2) &= \sigma T_1 \rightarrow \sigma T_2 \\ \sigma(x_1 : T_1, \dots, x_n : T_n) &= (x_1 : \sigma T_1, \dots, x_n : \sigma T_n) \end{aligned}$$

Two questions arise[2]:

1. “Are *all* substitution instances of t well typed?” That is, for every σ , do we have $\sigma\Gamma \vdash \sigma t : T$ for some T ?
2. “Is *some* substitution instance of t well typed?” That is, can we find a σ such that $\sigma\Gamma \vdash \sigma t : T$ for some T ?

The former is connected to parametric polymorphism, while the latter is related to type reconstruction problem. Usually one is interested in combination of the two, i.e. finding most general substitution that make term well-typed.

The main challenge here is how to combine together fragments that use unrelated type parameters into one term. This is typical situation with `let`-expression where some term is denoted by some name and then used in other term, possibly multiple times. This might require binding type parameters depending on the context where the term occurs. The following approach can be used to analyze this situation.

Term

Definition 2.2. Let Γ be a context and t a term. A *solution* for (Γ, t) is a pair (σ, T) such that $\sigma\Gamma \vdash \sigma t : T$.

For example, let $\Gamma = f:X, a:Y$ and $t = f a$. Then $([X \mapsto \text{Nat} \rightarrow \text{Nat}, Y \mapsto \text{Nat}], \text{Nat})$ is a solution for (Γ, t) .

The problem with this definition is that it is declarative and therefore quite impractical. To fix this alternative definition is used[2].

Definition 2.3. A *constraint set* C is a set of equations $\{S_i = T_i\}_{i \in \{1..n\}}$.

Definition 2.4. A substitution σ is said to *unify* an equation $S = T$ if the substitution instances σS and σT are identical. We say that σ *unifies* (or *satisfies*) C if it unifies every equation in C .

Definition 2.5. The *constraint typing relation* $\Gamma \vdash t : T \mid_X C$ is defined by the rules in next figure.

Informally, $\Gamma \vdash t : T \mid_X C$ can be read “term t has type T under assumptions Γ whenever constraints C are satisfied”.

The idea behind these rules is to make the term is typable independently of how its subterms are mixed together but balance it with the set of equations that will make sure that the different parts of the term agree. The rule CT-APP stands out for it introduces new type parameter. All rules require that there is no name capturing for type variables introduced.

Definition 2.6. Suppose that $\Gamma \vdash t : S \mid C$. A *solution* for (Γ, t, S, C) is a pair (σ, T) such that σ satisfies C and $\sigma S = T$.

The definition 2.6 provides us with additional information in convenient form. With its help it is much easier to obtain the answer whether the necessary type substitution for the given term exists (i.e. its type can be reconstructed).

But first the certain equivalence of definitions 2.2 and 2.6 should be established. This could be presented in the form of two theorems[2]:

Theorem 2.7 (Soundness of constraint typing). *Suppose that $\Gamma \vdash t : S \mid C$. If (σ, T) is a solution for (Γ, t, S, C) , then it is also a solution for (Γ, t) .*

$\frac{x:T \in \Gamma}{\Gamma \vdash x : T \quad _{\emptyset} \{ \}} \quad (\text{CT-VAR})$	$\frac{\Gamma \vdash t_1 : T \quad _{\mathcal{X}} C}{C' = C \cup \{T = \text{Nat}\}} \quad (\text{CT-PRED})$
$\frac{\Gamma, x:T_1 \vdash t_2 : T_2 \quad _{\mathcal{X}} C}{\Gamma \vdash \lambda x:T_1. t_2 : T_1 \rightarrow T_2 \quad _{\mathcal{X}} C} \quad (\text{CT-ABS})$	$\frac{\Gamma \vdash t_1 : T \quad _{\mathcal{X}} C}{C' = C \cup \{T = \text{Nat}\}} \quad (\text{CT-ISZERO})$
$\frac{\begin{array}{l} \Gamma \vdash t_1 : T_1 \quad _{\mathcal{X}_1} C_1 \quad \Gamma \vdash t_2 : T_2 \quad _{\mathcal{X}_2} C_2 \\ \mathcal{X}_1 \cap \mathcal{X}_2 = \mathcal{X}_1 \cap FV(T_2) = \mathcal{X}_2 \cap FV(T_1) = \emptyset \\ \mathbf{X} \notin \mathcal{X}_1, \mathcal{X}_2, T_1, T_2, C_1, C_2, \Gamma, t_1, \text{ or } t_2 \\ C' = C_1 \cup C_2 \cup \{T_1 = T_2 \rightarrow \mathbf{X}\} \end{array}}{\Gamma \vdash t_1 t_2 : \mathbf{X} \quad _{\mathcal{X}_1 \cup \mathcal{X}_2 \cup \{\mathbf{X}\}} C'} \quad (\text{CT-APP})$	$\Gamma \vdash \text{true} : \text{Bool} \quad _{\emptyset} \{ \} \quad (\text{CT-TRUE})$
$\Gamma \vdash 0 : \text{Nat} \quad _{\emptyset} \{ \} \quad (\text{CT-ZERO})$	$\Gamma \vdash \text{false} : \text{Bool} \quad _{\emptyset} \{ \} \quad (\text{CT-FALSE})$
$\frac{\Gamma \vdash t_1 : T \quad _{\mathcal{X}} C}{C' = C \cup \{T = \text{Nat}\}} \quad (\text{CT-SUCC})$	$\frac{\begin{array}{l} \Gamma \vdash t_1 : T_1 \quad _{\mathcal{X}_1} C_1 \\ \Gamma \vdash t_2 : T_2 \quad _{\mathcal{X}_2} C_2 \quad \Gamma \vdash t_3 : T_3 \quad _{\mathcal{X}_3} C_3 \\ \mathcal{X}_1, \mathcal{X}_2, \mathcal{X}_3 \text{ nonoverlapping} \\ C' = C_1 \cup C_2 \cup C_3 \cup \{T_1 = \text{Bool}, T_2 = T_3\} \end{array}}{\Gamma \vdash \text{if } t_1 \text{ then } t_2 \text{ else } t_3 : T_2 \quad _{\mathcal{X}_1 \cup \mathcal{X}_2 \cup \mathcal{X}_3} C'} \quad (\text{CT-IF})$

Figure 2.2: Constraint typing rules

Theorem 2.8 (Completeness of constraint typing). *Suppose that $\Gamma \vdash t : S \quad |_{\mathcal{X}} C$. If (σ, T) is a solution for (Γ, t) and $\text{dom}(\sigma) \cap \mathcal{X} = \emptyset$, then there is some solution (σ', T) for (Γ, t, S, C) such that $\sigma' \setminus \mathcal{X} = \sigma$.*

This means that if solution in one sense exists then there is solution in another sense, and either can be obtained from another.

Procedure that decides the constrain set is call *unification algorithm*. It is presented in the figure 2.3.

Unification algorithm has properties summarized in the next proposition.

Proposition 2.9. *The algorithm `unify` always terminates*

1. *`unify`(C) halts, either by failing or by returning a substitution, for all C ;*
2. *if `unify`(C) = σ , then σ is a unifier for C ;*
3. *if δ is a unifier for C , then `unify`(C) = σ with $\sigma \sqsubseteq \delta$.*

This reasoning can be adopted in the type system in the form of adjusted typing rules and more flexible in the for of algorithm. Rules are presented

$unify(C)$ = if $C = \emptyset$, then $[]$
 else let $\{S = T\} \cup C' = C$ in
 if $S = T$
 then $unify(C')$
 else if $S = X$ and $X \notin FV(T)$
 then $unify([X \mapsto T]C') \circ [X \mapsto T]$
 else if $T = X$ and $X \notin FV(S)$
 then $unify([X \mapsto S]C') \circ [X \mapsto S]$
 else if $S = S_1 \rightarrow S_2$ and $T = T_1 \rightarrow T_2$
 then $unify(C' \cup \{S_1 = T_1, S_2 = T_2\})$
 else
 fail

Figure 2.3: Unification algorithm

below

$$\frac{X \notin \chi \quad \Gamma, x : X \vdash t_1 : T \mid_{\chi} C}{\Gamma \vdash \lambda x. t_1 : X \rightarrow T \mid_{\chi \cup \{X\}} C} \text{ (CT-ABSINF)}$$

$$\frac{\Gamma \vdash [x \mapsto t_1]t_2 : T_2 \mid_{\chi} C}{\Gamma \vdash \mathbf{let} \ x = t_1 \ \mathbf{in} \ t_2 : T_2 \mid_{\chi} C} \text{ (CT-LETPOLY)}$$

Here the rule (CT-LETPOLY) is used instead of conventional

$$\frac{\Gamma \vdash [x \mapsto t_1]t_2 : T_2 \quad \Gamma \vdash t_1 : T_1}{\Gamma \vdash \mathbf{let} \ x = t_1 \ \mathbf{in} \ t_2 : T_2} \text{ (T-LETPOLY)}$$

It makes sure that each occurrence of t_1 would end up with individual copy of type variables. But there is a trade-off type checking complexity becomes exponential in the size of term if **let**-expression are nested in declaration part.

Rule-based type checking of expression **let** $x = t_1$ **in** t_2 can be replaced with the following algorithm [2]:

1. We use the constraint typing rules to calculate a type S_1 and a set C_1 of associated constraints for the right-hand side t_1 .
2. We use unification to find a most general solution σ to the constraints C_1 and apply σ to S_1 (and Γ) to obtain t_1 's principal type T_1 .

3. We generalize any variables remaining in T_1 . If $X_1 \dots X_n$ are the remaining variables, we write $\forall X_1 \dots X_n.T_1$ for the principal *type scheme*¹⁹ of t_1 .
4. We extend the context to record the type scheme $\forall X_1 \dots X_n.T_1$ for the bound variable x , and start type checking the body t_2 .
5. Each time we encounter an occurrence of x in t_2 , we look up its type scheme $\forall X_1 \dots X_n.T_1$. We now generate fresh type variables $Y_1 \dots Y_n$ and use them to instantiate the type scheme, yielding $[X_1 \mapsto Y_1, \dots, X_n \mapsto Y_n]T_1$, which we use as the type of x .

It is quite clear from the rule set of the system and presented algorithm that type schemes are only allowed in the head of `let`-expression. This means that for arguments of the lambda-abstraction only monotypes are allowed. This explains why earlier example produces an error.

The type system discussed in this subsection is known as Hindley–Milner or HM type system. It is implemented in ML family of programming languages. Modern programming language use HM system with slight modification for type inference, their type systems might be more expressive than HM system in certain aspects. For instance, OCaml includes algebraic data types. Haskell programming language implements typesystem for which HM system can be considered a restriction. It will be discussed in the next subsection.

2.3 First-class polymorphism

The next typesystem is called System F. It provides first-class polymorphism. Therefore it is called *polymorphic* lambda calculus²⁰. Its main feature is introduction of universal data types. Because of quantification of the type parameters it is also known as *second-order* lambda calculus²¹.

As usual we start with simpler system and extend it to new one. Additional type form appears:

$$T ::= X \quad | \quad \boxed{\forall X.T} \dots$$

There are two new operations on terms:

$$\begin{aligned} \boxed{\lambda X.t} & \quad \text{type abstraction} \\ \boxed{t [T]} & \quad \text{type application} \end{aligned}$$

Overall system is summarized in figure 2.4

¹⁹also called polytypes

²⁰Ukrainian equivalent: *поліморфне* лямбда-числення

²¹Ukrainian equivalent: *2-го порядку* лямбда-числення

$\rightarrow \forall$ Based on λ_{\rightarrow}

Syntax	Evaluation	$\boxed{t \longrightarrow t'}$
$t ::=$	<i>terms:</i>	
x	<i>variable</i>	
$\lambda x:T. t$	<i>abstraction</i>	$\frac{t_1 \longrightarrow t'_1}{t_1 t_2 \longrightarrow t'_1 t_2} \quad (\text{E-APP1})$
$t t$	<i>application</i>	$\frac{t_2 \longrightarrow t'_2}{v_1 t_2 \longrightarrow v_1 t'_2} \quad (\text{E-APP2})$
$\lambda X. t$	<i>type abstraction</i>	$(\lambda x:T_{11}. t_{12}) v_2 \longrightarrow [x \mapsto v_2]t_{12} \quad (\text{E-APPABS})$
$t [T]$	<i>type application</i>	$\frac{t_1 \longrightarrow t'_1}{t_1 [T_2] \longrightarrow t'_1 [T_2]} \quad (\text{E-TAPP})$
$v ::=$	<i>values:</i>	
$\lambda x:T. t$	<i>abstraction value</i>	$(\lambda X. t_{12}) [T_2] \longrightarrow [X \mapsto T_2]t_{12} \quad (\text{E-TAPPTABS})$
$\lambda X. t$	<i>type abstraction value</i>	
$T ::=$	<i>types:</i>	$\boxed{\Gamma \vdash t : T}$
X	<i>type variable</i>	
$T \rightarrow T$	<i>type of functions</i>	$\frac{x:T \in \Gamma}{\Gamma \vdash x : T} \quad (\text{T-VAR})$
$\forall X. T$	<i>universal type</i>	$\frac{\Gamma, x:T_1 \vdash t_2 : T_2}{\Gamma \vdash \lambda x:T_1. t_2 : T_1 \rightarrow T_2} \quad (\text{T-ABS})$
$\Gamma ::=$	<i>contexts:</i>	
\emptyset	<i>empty context</i>	$\frac{\Gamma \vdash t_1 : T_{11} \rightarrow T_{12} \quad \Gamma \vdash t_2 : T_{11}}{\Gamma \vdash t_1 t_2 : T_{12}} \quad (\text{T-APP})$
$\Gamma, x:T$	<i>term variable binding</i>	$\frac{\Gamma, X \vdash t_2 : T_2}{\Gamma \vdash \lambda X. t_2 : \forall X. T_2} \quad (\text{T-TABS})$
Γ, X	<i>type variable binding</i>	$\frac{\Gamma \vdash t_1 : \forall X. T_{12}}{\Gamma \vdash t_1 [T_2] : [X \mapsto T_2]T_{12}} \quad (\text{T-TAPP})$

Figure 2.4: System F

Typing rule (T-TAbs) introduces quantification of type variables. Evaluation rules are basically intact. Type safety theorems can be proven[2]. They retain the same form as for system λ_{\rightarrow} .

Theorem 2.10 (preservation). *If $\Gamma \vdash t : T$ and $t \rightarrow t'$, then $\Gamma \vdash t' : T$*

Theorem 2.11 (progress). *If t is a closed, well-typed term (that is, $\vdash t : T$ for some T). Then either t is a value or else there is some t' with $t \rightarrow t'$.*

System F can also be shown to have a normalization property, i.e. that typable terms are normalizing.

Now we are able to fully explain the earlier example regarding first-class polymorphism. Construction used can be presented in system F with the following well-typed term:

$$\begin{aligned} \text{id} &\equiv \Lambda t. \lambda x : t. x \\ f &\equiv \lambda g : \forall t. t \rightarrow t. (g \text{ [int] } 1, g \text{ [real] } 1.0) \end{aligned}$$

This can be implemented as the following Haskell code

```
f :: (forall t. t -> t) -> (Integer, Double)
f g = (g 1 :: Integer), g 1.0
res = f id
```

Type reconstruction problem can be formulated for system F. The erasure for system F is defined as follows:

$$\begin{aligned} \text{erase}(x) &= x \\ \text{erase}(\lambda x : T. t) &= \lambda x. \text{erase}(t) \\ \text{erase}(t_1 t_2) &= \text{erase}(t_1) \text{erase}(t_2) \\ \text{erase}(\lambda X. t_2) &= \text{erase}(t_2) \\ \text{erase}(t_1 [T_2]) &= \text{erase}(t_1) \end{aligned}$$

Decidability for type reconstruction problem for system F was one of long standing problems in type theory. Finally the following was shown:

Theorem 2.12 (Well, 1994). *t is undecidable whether, given a closed term m of the untyped lambda-calculus, there is some well-typed term t in System F such that erase(t) = m.*

Therefore, unfortunately, in general case type annotations have to be provided for terms in System F.

Given a type with quantifiers consider the path from the root to quantifier when the type is drawn as a tree.

Definition 2.13. The *rank* of a quantifier in the type is the number of passes to the left of arrows in this path.

For example, in type term $\forall T. (\forall X. X) \rightarrow T$ first quantifier has rank 0 and second has rank 1.

Definition 2.14. Type is called of *rank* $k+1$ if its greatest rank of quantifiers is k .

Definition 2.15. Polymorphism is called *rank- n* ²², where n is a number, if it allows at most rank- n types. Rank-1 polymorphism is also called *prenex*. Polymorphism is called *arbitrary-rank (higher-rank)* if it allows any System F types, that is rank- n types, for any n .

Prenex polymorphism is called so because all type terms can be written down in prenex form²³. That is exactly the form used in `let`-polymorphism for type schemes or polytypes.

Definition 2.16. Polymorphism is called *predicative*²⁴ if type variables range only over quantifier-free types, i.e. monotypes. Polymorphism is called *impredicative*²⁵ if type variables are allowed to range over quantified types including the very type in which they occur.

²²Ukrainian equivalent: *рангу n*

²³when all quantifiers are in front

²⁴Ukrainian equivalent: *предикативний*

²⁵Ukrainian equivalent: *непредикативний*

Chapter 3

Selected type-based programming techniques

3.1 Elements of generic programming

Generic programming is a style of programming that strives to present algorithm and data structures in most abstract form without losing the efficiency. The central notion of generic programming is generic algorithm, which is procedural schemata parameterized by data types and independent of their specific representation[11]. This way it is possible to group solutions of the similar problems compactly and effectively.

Key ideas of generic programming are following.

- Expressing algorithms with minimal assumptions about data abstractions, and vice versa, thus making them as interoperable as possible
- Lifting of a concrete algorithm to as general a level as possible without losing efficiency; i.e., the most abstract form such that when specialized back to the concrete case the result is just as efficient as the original algorithm
- When the result of lifting is not general enough to cover all uses of an algorithm, additionally providing a more general form, but ensuring that the most efficient specialized form is automatically chosen when applicable.
- Providing more than one generic algorithm for the same purpose and at the same level of abstraction, when none dominates the others in efficiency for all inputs. This introduces the necessity to provide suffi-

ciently precise characterizations of the domain for which each algorithm is the most efficient.

3.1.1 Basic notions of generic programming

Like many other paradigms generic programming could be viewed at the following levels: general ideas, principle and notions of it, mechanisms in particular programming language that facilitate programming in this style and specific systems (libraries) that are built according to it.

The notion of *concepts*¹ is introduced. Typically they are compositions of requirements that should be satisfied by type parameters used in generic algorithm declaration. It was mentioned earlier that requirements could be syntactic and semantic. If syntactic requirements are met by types it is said that the concept is *satisfied*. If additionally semantic requirements are met then the concept is *modeled*. Logic of the algorithm relies on concepts used. Concepts can be unary, but can involve multiple type parameters as well. If extension of concept C' is included in extension of concept C then C' is said to *refine*² concept C . Concepts organized by refinement relation form the taxonomy. Authors of generic programming draw parallels between generic programming and mathematical construction like algebras and algebraic structures [12, 13].

Means through which programming language provides support for generic programming constructions are called *generics*³ [10]. Usually they realize some kind of parametric polymorphism.

3.1.2 Example of algorithm generalization

One usually doesn't start with generic version. Generic algorithms are obtained from concrete cases by gradually lifting level of abstraction. In this process concepts are taking shape. Sometimes wrong generalization would be made, then extending algorithm to another case would lead to adjusted version.

For many reasons C++ became the de facto standard programming language for generic programming. Therefore code samples in this section would be presented in C++.

Let us consider two following normal function

```
int sum(int* array, int n) {
```

¹Ukrainian equivalent: *концепти*

²Ukrainian equivalent: *уточнює*

³Ukrainian equivalent: *дженерики*

```

    int result = 0;
    for (int i = 0; i < n; ++i)
        result = result + array[i];
    return result;
}

float sum(float* array, int n) {
    float result = 0;
    for (int i = 0; i < n; ++i)
        result = result + array[i];
    return result;
}

```

Though the work they do is for different types they are very similar. We can generalize them to arguments of numeric types:

```

template<typename T>
T sum(T* array, int n) {
    T result = 0;
    for (int i = 0; i < n; ++i)
        result = result + array[i];
    return result;
}

```

This version does not take into account following case

```

std::string concatenate(std::string* array, int n) {
    std::string result = "";
    for (int i = 0; i < n; ++i)
        result = result + array[i];
    return result;
}

```

In order to do it variable initialization should be corrected:

```

template<typename T>
T sum(T* array, int n) {
    T result = T();
    for (int i = 0; i < n; ++i)
        result = result + array[i];
    return result;
}

```

Now we want to allow other types of collections:

```

template<typename Container, typename T>
T sum(const Container& array, int n) {
    T result = T();
}

```

```

    for (int i = 0; i < n; ++i)
        result = result + array[i];
    return result;
}

```

This generalization was in wrong direction for it is not suitable for next case (lists)

```

template<typename T>
struct node {
    node<T>* next;
    T data;
};

template<typename T>
T sum(node<T> *first, node<T> *last) {
    int s = T();
    for (; first != last; first = first->next)
        s = s + first->data;
    return s;
}

```

Here we see that use of array index operator in function template contradicts to linked list approach because list elements are not randomly accessed. In order to cover this case we need to reconsider they way how the collection is traversed for earlier cases. We come to the following function

```

template<typename T>
T sum(T* array, int n) {
    T result = T();
    for (T* current = array; current != array + n;
        ++current)
        result = result + *current;
    return result;
}

```

Then the generalization of iteration through elements would shape as follows⁴

```

template<typename I, typename T>
T sum(I start, I end, T init) {
    for (I current = start; current != end;
        current = next(current))
        init = init + get(current);
    return init;
}

```

⁴actually in C++ dereference and increment operators could retain

```
}
```

During specialization (application of algorithm to lists) compiler would be looking for the following definitions:

```
template<typename T>
struct node {
    node<T>* next;
    T data;
};

template<typename T>
node<T>* next(node<T>* n) { return n->next; }

template<typename T>
T get(node<T>* n) { return n->data; }
```

They provide adaptation of general schemata to the linked-lists.

Based on the last version of algorithm we derive several requirements to the type parameters

- *T must have an additive operator +*
- *T must have an assignment operator*
- *I must have an inequality operator !=*
- *I must have a copy constructor*
- *I must have an operation next() that moves to the next value in the sequence*
- *I must have an operation get() that returns the current value (of type T).*

For these requirements properly named concepts could be introduced. They are given in table 3.1. Then combining information for this generic algorithms the iterator concept can be formed. It is presented in table 3.2. Then we consider another generic algorithm

```
// Requirements: Iterator<I, T>
template<typename I, typename T>
int distance(I start, I end) {
    int i = 0;
    for (; start != end; ++start)
        ++i;
    return i;
}
```

Concept	Requirement
CopyConstructible<T>	T must have a copy constructor.
Assignable<T>	T must have an assignment operator.
Addable<T>	T must have an additive operator +.
EqualityComparable<T>	T must have operators ==, != comparing two Ts and returning a bool.

Table 3.1: Named concepts

Concept	Requirement
	EqualityComparable<I> CopyConstructible<I>
Iterator<I, T>	Assignable<I> I has operation next() that moves to the next value in the sequence I has operation get() that returns the current value of T

Table 3.2: First version of iterator concept

And it is obvious that second type parameter is redundant for it. Therefore we add requirement for associated to obtain update algorithm.

```
template<typename I>
typename iterator_traits<I>::value_type
sum(I start, I end, typename iterator_traits<I>::value_type init) {
    for (I current = start; current != end; current = next(current))
        init = init + get(current);
    return init;
}
```

For language supporting concepts algorithm would be similar to the next.

```
template<InputIterator Iter>
```

Concept	Requirement
	EqualityComparable<I> CopyConstructible<I>
Iterator<I>	Assignable<I> value_type is an associated type next, get returns value_type

Table 3.3: Updated version of iterator concept

```

        where Addable<Iter::value_type>
              && Assignable<Iter::value_type>
Iter::value_type
sum(Iter first, Iter last, Iter::value_type init) {
    for (; first != last; first = next(first))
        init = init + get(first);
    return init;
}

```

Finally we can write algorithm generalized also by binary operation to obtain.

```

template<InputIterator Iter,
        BinaryOperation<Iter::value_type, Iter::value_type> Op>
    where Assignable<Iter::value_type, Op::result_type>
Iter::value_type
accumulate(Iter first, Iter last, Iter::value_type init, Op op) {
    for (; first != last; first = next(first))
        init = op(init, get(first));
    return init;
}

```

Up to C++20 version of the language this features were purely experimental, therefore unstandardized and used only by geeks. Concepts were included in C++20 standard. They are defined in terms of constant expressions. Their syntax is similar to presented but slightly different. This concerns constrained template parameters, namely how n-ary concepts are referred to. Naming conventions used for them are different also. New feature in standard library use concepts heavily.

3.1.2.1 Generic programming in STL

Projects like standard template library (STL) were testing ground for generic programming. STL library is still one of the base case-studies of generic programming techniques. Many fruitful ideas and generic algorithm are implemented in this system. Concise account of it from generic programming perspective is presented in [10].

3.2 Datatype-generic programming

In generic programming we try to achieve more possibilities for parametrization without expanding possibilities for uncaught errors. *Datatype-generic programming*⁵ (DGP), on the other hand, is a special approach to genericity

⁵Ukrainian equivalent: *програмування параметризоване за класами алгебр*

based on parametrization by *shape* (polytypism). DGP is quite similar to generic programming techniques. The prefix ‘Data’ is used to highlight the differences from them. Extended account on DGP can be found in [14] and [15].

3.2.1 Roy–Floyd–Warshall algorithm

Let us start with motivational example.

Reachability problem In the following code solution for reachability in the graph problem is considered. It is known as Warshall’s algorithm. It should be mentioned that it find routes for all pairs of vertices. It starts from adjacency matrix of initial graph and finishes with same matrix for its transitive closure.

```
for each  $k, 1 \leq k \leq N$ 
do for each pair  $(i, j), 1 \leq i, j \leq N$ 
  do  $a_{ij} := a_{ij} \vee (a_{ik} \wedge a_{kj})$ 
  end_for
end_for
```

In order to talk about correctness of this algorithm its pre- and post-conditions should be written down. Without this information algorithm definition is not fully meaningful. Therefore, precondition:

- a matrix $N \times N$ of Bool
- $a_{ij} = \text{true}$ if there an edge (i, j) .

Postcond:

- $a_{ij} = \text{true}$ if there path $[i, j]$.

Least-cost path problem Similar to the previous algorithm this one looks for shortest path for all pairs of vertices. It is known as Floyd’s algorithm.

```
for each  $k, 1 \leq k \leq N$ 
do for each pair  $(i, j), 1 \leq i, j \leq N$ 
  do  $a_{ij} := a_{ij} \downarrow (a_{ik} \oplus a_{kj})$ 
  end_for
end_for
```

where $x \downarrow y$ is a minimum of x and y . Once again following information is provided for completeness. Precond:

- a matrix $N \times N$ of $\text{Int} \cup \{\infty\}$.
- a_{ij} is a cost passing edge (i, j) .
- $a_{ij} = \infty$ if there is no edge.

Postcond:

- a_{ij} is the least cost passing from i to j .

Bottleneck problem The last example is maximin paths. The solution is presented below

```

for each  $k, 1 \leq k \leq N$ 
do for each pair  $(i, j), 1 \leq i, j \leq N$ 
  do  $a_{ij} := a_{ij} \uparrow (a_{ik} \downarrow a_{kj})$ 
  end_for
end_for

```

where additionally $x \uparrow y$ is a maximum of x and y . Precond:

- a matrix $N \times N$ of $\text{Int} \cup \{\infty\}$.
- a_{ij} is a least height of bridge on edge (i, j) .
- $a_{ij} = 0$ if there is no edge.

Postcond:

- a_{ij} is the max possible height of vehicle to pass from i to j .

Generic path-algorithm (transitive closure algorithm) All three algorithms have similar shapes. The only difference are initial values and the concrete binary operations used in each of them. If other algorithms are also taken into account we achieve abstract (general) form of the algorithm.

```

for each  $k, 1 \leq k \leq N$ 
do for each pair  $(i, j), 1 \leq i, j \leq N$ 
  do  $a_{ij} := a_{ij} \oplus (a_{ik} \otimes a_{kk}^* \otimes a_{kj})$ 
  end_for
end_for

```

It is possible to talk about applicability (executability) and correctness of the algorithm. Assume that the type of the matrix elements is A . Then the algorithm applies when this type supports operations used in algorithm. This situation is studied in abstract algebra. The list of operations supported

and their *arities*⁶ are called signature. Hence for applicability type of the matrix elements A is a carrier of algebra with signature $(\oplus, \otimes, *)$. Once again we specify pre- and post-conditions relying to some extent on mentioned operations. Precond:

- a matrix $N \times N$ of A .
- $a_{ij} = \sum\{e : e \text{ is an edge from } i \text{ to } j : \text{label } e\}$.

Postcond:

- $a_{ij} = \sum\{p : p \text{ is path from } i \text{ to } j : \text{weight } p\}$
 $\text{weight } [] = 1$
 $\text{weight}(e : p) = \text{label } e \otimes \text{weight } p$

Now for the algorithm to make sense certain requirements to algebraic operations should be met. Algorithm is correct if $(A, \oplus, \otimes, *, 0, 1)$ is a *regular algebra*⁷, which means following.

- $(0, 1, \oplus, \otimes)$ form semiring with unit structure;
- \oplus should be idempotent
 - partial order on A : $a \leq b \Leftrightarrow a \oplus b = b$;
- $*$ – least fixed-point operator for \otimes :

$$\begin{aligned} 1 \oplus a \otimes a^* &\leq a^* \geq a^* \otimes a \oplus 1 \\ a \otimes x &\leq x \rightarrow a^* \otimes x \leq x \\ x \otimes a &\leq x \rightarrow x \otimes a^* \leq x \end{aligned}$$

In abstract algebra these requirements are usually called axioms. Some information on this matter is provided in the next subsection.

3.2.2 Algebras & Morphisms

Good account of algebras is provided in [16]. Here only the basic notions are recalled. (Concrete) Algebra is a set A (called) with a number of operations on it. Operations f is a function of type $A^k \rightarrow A$, k is operation arity. There

⁶Ukrainian equivalent: *арності*

⁷Ukrainian equivalent: *регулярна алгебра*

are prefix, infix, suffix forms of notation. In mathematics and programming binary operations are usually written infix.

Here are a few algebra examples⁸. Most of them are self-explaining. Symbol ++ denotes word concatenation.

- $(N, 0, (+))$ with $0 :: 1 \rightarrow N$, $(+) :: N \times N \rightarrow N$
- $(N, 0, (\uparrow))$ with $0 :: 1 \rightarrow N$, $(\uparrow) :: N \times N \rightarrow N$
- $(R, 1, (\times))$ with $1 :: 1 \rightarrow R$, $(\times) :: R \times R \rightarrow R$
- $(B, \top, (\equiv))$ with $\top :: 1 \rightarrow B$, $(\equiv) :: B \times B \rightarrow B$
- $(B, \perp, (\vee))$ with $\perp :: 1 \rightarrow B$, $(\vee) :: B \times B \rightarrow B$
- $(B, \top, (\wedge))$ with $\top :: 1 \rightarrow B$, $(\wedge) :: B \times B \rightarrow B$
- $(A^*, \varepsilon, (\text{++}))$ with $\varepsilon :: 1 \rightarrow A^*$, $(\text{++}) :: A^* \times A^* \rightarrow A^*$

Datatypes in Haskell One of the features of the programming languages is the possibility to define self referential types. Therefore such types are called *inductive* types⁹. Algebras arise in Haskell naturally from type definition using `data` keyword. Consider the following example in Haskell.

```
data Nat = zero | succ Nat
```

First word of alternatives defined in the type are called value constructors¹⁰ The following conversation shows their type:

```
*Main> :t succ
succ :: Nat -> Nat
```

This is exactly the type of unary operation in the algebra with the carrier `Nat`. This way data type and value constructors form an algebra

- $(\text{Nat}, \text{zero}, \text{succ})$ with $\text{zero} :: 1 \rightarrow \text{Nat}$, $(\text{succ}) :: \text{Nat} \rightarrow \text{Nat}$

Therefore such types are also called *algebraic*¹¹.

⁸by Haskell convention operators are parenthesized

⁹Ukrainian equivalent: *індуктивні типи*

¹⁰For illustration their first letters are lowercased

¹¹Ukrainian equivalent: *алгебраїчні*

Algebra morphisms Algebra *morphisms*¹² are the algebraic structure-preserving maps between two algebras. They depend on the class of algebras. For each class of algebras we have specific type of morphisms. Both source and target algebra should belong to the same class.

Here is an example. The algebra morphisms $\text{exp} : (N, 0, (+)) \rightarrow (R, 1, (\times))$ is defined in the following way (exponentiation):

$$\begin{aligned}\text{exp } 0 &= 1 \\ \text{exp}(x + y) &= (\text{exp } x) \times (\text{exp } y)\end{aligned}$$

Another example is the function $\text{length} : (A^*, \varepsilon, (++)) \rightarrow (N, 0, (+))$ between strings and numbers

$$\begin{aligned}\text{length } \varepsilon &= 0 \\ \text{length}(s ++ q) &= \text{length } s + \text{length } q\end{aligned}$$

Defining function on inductive datatypes Functions can be defined on inductive types. For this *pattern matching*¹³ is used where for each kind of values of the type certain processing specified. Consider example in Haskell:

```
data Nat = zero | succ Nat
— pattern matching
pred zero = zero
pred (succ x) = x
even zero = True
even (succ x) = not (even x)
```

This way morphisms between algebras can be constructed.

Algebras of some class with their morphisms form a category. Another level of abstraction can be form using this categories.

Types parameterization also fits this scheme. Consider the notion of type constructors in Haskell. These entities:

- create new types from existing;
- thus have type parameters.

The `List` defined below gives an example of type constructor.

```
data List a = nil | cons a (List a)
```

Action of constructor `List` can be extended to functions

¹²also homomorphisms

¹³Ukrainian equivalent: *зіставлення зі зразком*

```
mapList f = h where
  h nil = nil
  h (cons x xs) = cons (f x) (h xs)
```

3.2.3 Type constructors and functors

Type constructors like `List` enjoy good property. Name they preserve that categorical structure.

In order to talk about it we need to introduce couple notions first.

Definition 3.1. A *category* C consists of

- collection of *objects* $\text{Obj } C$
- collection of *morphisms* (*arrows*) $\text{Mor } C$

such that

- each $f \in \text{Mor } C$ is assigned a type $f: a \rightarrow b$, where $a, b \in \text{Obj}$
- for $f: a \rightarrow b$, $g: b \rightarrow c$ the *composition* is defined $f \circ g: a \rightarrow c$
- for all $f: a \rightarrow b$, $g: b \rightarrow c$; $h: c \rightarrow d$

$$f \circ (g \circ h) = (f \circ g) \circ h$$

- for each $a \in \text{Obj } C$ there is the *identity* $1_a: a \rightarrow a$ such that for $f: a \rightarrow b$

$$f \circ 1_b = f = 1_a \circ f$$

Type constructors generalize to the notion of functor. A functor $F: C \rightarrow D$ is a map between categories that preserves structures

$$\begin{aligned} F: \text{Obj } C &\rightarrow \text{Obj } D \\ F: C(a, b) &\rightarrow D(Fa, Fb) \end{aligned}$$

such that

$$\begin{aligned} F1_a &= 1_{Fa} \\ F(f \circ g) &= Ff \circ Fg. \end{aligned}$$

n -ary functors Using product categories functors can be further generalized to n -ary functors. Functor $F: C_1 \times \dots \times C_n \rightarrow D$ is a map

$$F: \prod_{i=1}^n \text{Obj } C_i \rightarrow \text{Obj } D$$

$$F: \prod_{i=1}^n C_i(a_i, b_i) \rightarrow D(Fa_1 \dots a_n, Fb_1 \dots b_n)$$

such that

$$F1_{a_1} \dots 1_{a_n} = 1_{Fa_1 \dots a_n}$$

$$F(f_1 \circ g_1 \dots f_n \circ g_n) = Ff_1 \dots f_n \circ Fg_1 \dots g_n.$$

In Haskell programming language there is the category \mathcal{H}_{ASK} . Object and morphisms are Haskell types¹⁴ and functions respectively:

$$\text{Obj}_{\mathcal{H}_{ASK}} = \{T \mid T \text{ - is a type}\}$$

$$\mathcal{H}_{ASK}(T, S) = \{f \mid f \text{ - is a function of type } T \rightarrow S\}$$

Composition defined as usual function composition:

$f \circ g = h$ where $h \ x = f \ (g \ x)$

$$g \circ f = f \circ g$$

3.2.4 Polynomial Functors

Polynomial functors¹⁵ play important role because most algebraic data type definitions in Haskell can be derived using them. The reason for this is that basic building blocks like alternatives and multiple parameters in constructors translate to well-known categorical notions of sum and product. And also constants (nullary data constructor) obviously translate to constant functors. Finally, flexibility with type parameters is achieved using extraction functors.

3.2.4.1 Const and extraction functors

Let $a \in \text{Obj } C$. Define constant functor a^K as follows.

$$a^K: x_1 \dots x_n \mapsto a$$

$$a^K: f_1 \dots f_n \mapsto 1_a$$

Then $a^K: C^n \rightarrow C$ is an n -functor. Here is an example of code for the case when $n = k = 2$.

¹⁴due to limitations some subset of types is used

¹⁵Ukrainian equivalent: *поліноміальний* функтори

```
data Ex22 a b = In22 b
map22 f g (In22 x) = In22 (g x)
```

Define the family of extraction functors $\text{Ex}_k^n: C_1 \times \dots \times C_n \rightarrow C_k$ such that

$$\begin{aligned}\text{Ex}_k^n: x_1 \dots x_n &\mapsto x_k \\ \text{Ex}_k^n: f_1 \dots f_n &\mapsto f_k\end{aligned}$$

Particularly $I_C = \text{Ex}_1^1: C \rightarrow C$ (identity functor):

$$\begin{aligned}I_C x &= x \\ I_C f &= f\end{aligned}$$

3.2.4.2 Sum functor

Sum functor captures the notion of alternatives. Consider the following fragment in Haskell¹⁶.

```
data a + b = inl a | inr b
```

```
f + g = h where
  h (inl x) = inl (f x)
  h (inr x) = inr (g x)
```

```
f ∇ g = h where
  h (inl x) = f x
  h (inr x) = g x
```

This can be represented in, so called, pointfree form or style:

$$\begin{aligned}h = f + g &\Leftrightarrow \begin{cases} \text{inl} \circledast h = f \circledast \text{inl} \\ \text{inr} \circledast h = g \circledast \text{inr} \end{cases} \\ h = f \nabla g &\Leftrightarrow \begin{cases} \text{inl} \circledast h = f \\ \text{inr} \circledast h = g \end{cases}\end{aligned}$$

This can be understood as the definition of the function as the unique solution of the corresponding equation. In category theory construction with these properties is called coproduct. Specifically in category theory respective equations are typically represented as commutative diagrams.

¹⁶To define infix type constructors `TypeOperators` option is required here

This can be used to prove identities by calculating. If we have to prove something for $f \nabla g \circ h$. Let us just try to express it as “junk” itself:

$$f \nabla g \circ h = \alpha \nabla \beta \Leftrightarrow \begin{cases} \text{inl} \circ f \nabla g \circ h = \alpha \\ \text{inr} \circ f \nabla g \circ h = \beta \end{cases} \Leftrightarrow \begin{cases} \alpha = f \circ h \\ \beta = g \circ h \end{cases}$$

That is $f \nabla g \circ h = (f \circ h) \nabla (g \circ h)$

3.2.4.3 Product functor

For type products representation we can use either tuples or infix variant of type constructor. Here is the code fragment of the former.

```
type a × b = (a, b)

exl :: a × b → a
exl (x, y) = x
exr (x, y) = y
(f × g) (x, y) = (f x, g y)
(f Δ g) x = (f x, g x)
```

The last two identities can be rewritten in pointfree form as follows.

$$h = f \times g \Leftrightarrow \begin{cases} h \circ \text{exl} = \text{exl} \circ f \\ h \circ \text{exr} = \text{exr} \circ g \end{cases}$$

$$h = f \nabla g \Leftrightarrow \begin{cases} h \circ \text{exl} = f \\ h \circ \text{exr} = g \end{cases}$$

3.2.4.4 Combining functors

For

$$F_i: C_1 \times \cdots \times C_n \rightarrow D_i, \quad i = \overline{1, m}$$

$$G: D_1 \times \cdots \times D_m \rightarrow E$$

there is functor $H: C_1 \times \cdots \times C_n \rightarrow E$ such that

$$H a_1 \dots a_n = G (F_1 a_1 \dots a_n) \dots (F_m a_1 \dots a_n)$$

$$H f_1 \dots f_n = G (F_1 f_1 \dots f_n) \dots (F_m f_1 \dots f_n)$$

i.e.

$$H_{-1} \dots_{-n} = G (F_{1-1} \dots_{-n}) \dots (F_{n-1} \dots_{-n})$$

3.2.4.5 Polynomial functor definition and example

Definition 3.2. A *polynomial functor* is a functor obtained by combination of extract and constant functors.

EXAMPLE. Let's get back to the list example

```
data List a = nil | cons a (List a)
```

Its *pattern* functor¹⁷ is the binary polynomial functor L defined on objects as follows.

$$L a z = 1 + a \times z$$

with obvious definition on morphisms.

3.2.5 Functor algebras

Let F be an endofunctor on category C , i.e. $F: C \rightarrow C$.

Definition 3.3. An F -algebra is a pair (a, φ) , such that $\varphi: Fa \rightarrow a$.

Definition 3.4. An F -algebra morphism $f: (a, \varphi) \rightarrow (b, \psi)$ is a morphism $f: a \rightarrow b$ of category C , such that $f \circ \varphi = \psi \circ Ff$. Equivalently following diagram

$$\begin{array}{ccc} Fa & \xrightarrow{Ff} & Fb \\ \varphi \downarrow & & \downarrow \psi \\ a & \xrightarrow{f} & b \end{array}$$

commutes.

We can compose F -algebra morphisms using morphism composition in C .

Proposition 3.5. F -algebra with F -algebra morphisms form a category. Let's denote it $(F: C)$.

EXAMPLE. Let $\text{Square} : \mathbf{Set} \rightarrow \mathbf{Set}$ be a doubling functor $x \mapsto x \times x$. Then $(\text{Square}, \mathbf{Set})$ is a category of magmas, i.e. (unconstrained) algebras with single binary operation.

EXAMPLE. In general, if $F : \mathbf{Set} \rightarrow \mathbf{Set}$ is a polynomial functor then (F, \mathbf{Set}) is a category of universal algebras with signature determined by functor F .

¹⁷Ukrainian equivalent: *каркасний* функтор

Definition 3.6. An initial object of category $(F : C)$ is called *initial F -algebra*, denoted (μ_F, in_F) .

Definition 3.7. Unique morphism from initial F -algebra to any F -algebra (a, ϕ) is called *catamorphism*, denoted $(\downarrow\phi)$.

EXAMPLE. Consider algebra with unary operation f and constant e . To it corresponds the polynomial functor $F: x \mapsto 1 + x$. The initial algebra in this case is $(N, 1 \nabla s)$. This can be demonstrated by commutative diagram

$$\begin{array}{ccc} N + 1 & \xrightarrow{h+1} & X + 1 \\ \downarrow s \nabla 1 & & \downarrow f \nabla a \\ N & \xrightarrow{h} & X \end{array}$$

which is the categorical reformulation of recursive definition of h .

Lemma 3.8 (Plotkin). *An initial F algebra is an isomorphism.*

Proof. Let (μ_F, in_F) be an initial F -algebra. Consider the morphism $F(\text{in}_F)$. Clear it is an F -algebra $(F\mu_F, F(\text{in}_F))$. We denote by $\text{out}_F: F\mu_F \rightarrow \mu_F$ a catamorphism for this algebra. Let's prove that out_F is an inverse of in_F . Consider following commutative diagram

$$\begin{array}{ccccc} F\mu_F & \xrightarrow{F\text{out}_F} & FF\mu_F & \xrightarrow{F\text{in}_F} & F\mu_F \\ \downarrow \text{in}_F & & \downarrow F\text{in}_F & & \downarrow \text{in}_F \\ \mu_F & \xrightarrow{\text{out}_F} & F\mu_F & \xrightarrow{\text{in}_F} & \mu_F \end{array} .$$

Outer square defines catamorphism $(\downarrow\text{in}_F)$. By universal property $(\downarrow\text{in}_F) = \text{id}_{\mu_F}$. Then immediately from left square:

$$\text{out}_F \circ \text{in}_F = F(\text{in}_F \circ \text{out}_F) = \text{id}_{F\mu_F}$$

□

This construction has remarkable property. If the functor F is parameterized by additional parameters, i.e. it is an n -functor, and it satisfies certain conditions then initial F -algebra μ_F is also functorial. Many inductive types (e.g. in Haskell) can be formalized this way. This is the case for types

like lists. It also explains uniqueness of recursive definitions on inductive datatypes.

The dual construction to initial algebras are final co-algebras, final objects of category of co-algebras. Using them coinductive datatypes such as streams are formalized.

Bibliography

- [1] Luca Cardelli and Peter Wegner. On understanding types, data abstraction, and polymorphism. *ACM Computing Surveys*, 17(4):471–522, 1985.
- [2] Benjamin C. Pierce. *Types and Programming Languages*. MIT Press, 2002.
- [3] Luca Cardelli. Type systems. In Allen B. Tucker, editor, *The Computer Science and Engineering Handbook*, chapter 97, pages 2277–2308. CRC Press, 2004.
- [4] Niklaus Wirth. A brief history of software engineering. *IEEE Annals of the History of Computing*, 30(3):32–39, 2008.
- [5] R. M. Burstall. Proving properties of programs by structural induction. *Computer Journal*, 12(1):41–48, 1969.
- [6] Kevin Klement. Russell’s 1903–05 anticipation of the lambda calculus. *History and Philosophy of Logic*, 24:15–37, 2003.
- [7] Bernard Linsky. The notation in principia mathematica. In Edward N. Zalta, editor, *Stanford Encyclopedia of Philosophy*. 2011.
- [8] Christopher Strachey. Fundamental concepts in programming languages. *Higher-Order and Symbolic Computation*, 13(1/2):11–49, 2000. Reprinted from lecture notes dated 1968, with a foreword [17].
- [9] Kim Bruce, Luca Cardelli, Giuseppe Castagna, The Hopkins Object Group, Gary T. Leavens, and Benjamin Pierce. On binary methods. *Theory and Practice of Object Systems*, 1(3):221–242, 1995.
- [10] Jeremy G. Siek. *A Language for Generic Programming*. PhD thesis, Indiana University, August 2005.

- [11] David Musser and Alexander A. Stepanov. Generic programming. In *Symbolic and algebraic computation: ISSAC '88*, pages 13–25. Springer, 1988.
- [12] Alexander Stepanov and Paul McJones. *Elements of Programming*. Addison-Wesley Professional, 1st edition, 2009.
- [13] Alexander A. Stepanov and Daniel E. Rose. *From Mathematics to Generic Programming*. Addison-Wesley Professional, 1st edition, 2014.
- [14] Roland Carl Backhouse, Patrik Jansson, Johan Jeuring, and Lambert Meertens. Generic programming: An introduction. In *3rd International Summer School on Advanced Functional Programming*, volume 1608 of *LNCS*, pages 28–115. Springer-Verlag, 1999.
- [15] Jeremy Gibbons. Datatype-generic programming. In *Proceedings of the 2006 International Conference on Datatype-generic Programming*, volume 4719 of *LNCS*, pages 1–71, Berlin, Heidelberg, 2007. Springer-Verlag.
- [16] D. Sannella and A. Tarlecki. *Foundations of Algebraic Specification and Formal Software Development*. Monographs in Theoretical Computer Science. An EATCS Series. Springer-Verlag Berlin Heidelberg, 2012.
- [17] Peter D. Mosses. A foreword to ‘Fundamental concepts in programming languages’. *Higher-Order and Symbolic Computation*, 13(1/2):7–9, 2000.

Glossary

abstract syntax tree абстрактне синтаксичне дерево. 7, 11, *see* parse tree

algebra алгебра.

carrier носій. 52

regular регулярна. 52

arity арність. 52

coercion неявне приведення типів. 30

concept концепт. 44, 49

contravariant контраваріантність. 33

covariant коваріантність. 33

datatype-generic programming програмування параметризоване за класами алгебр. 49

derivation виведення. 14

evaluation виведення (кроку) обчислення. 14

induction on індукція за виведеннями. 14

tree дерево. 14

typing виведення типу. 18

evaluation обчислення.

eager жадібне. 23

greedy жадібне. 23

lazy лінійні. 23

one-step один крок. 13, 14

strategy стратегія. 14, 22

strict строге. 23

first-class citizen об'єкт першого класу. 34

functor функтор.

pattern каркасний. 59

polynomial поліноміальний. 56

generic programming параметризоване програмування. 34, 43
generics дженерики. 44

induction індукція.
 on derivations за виведеннями. 14

lambda calculus лямбда-числення.
 polymorphic поліморфне. 39
 second-order 2-го порядку. 39
 simply-typed з простою типізацією. 24

normal form нормальна форма. 15, *see term*

overloading довизначення, перевантаження. 30

pattern matching зіставлення зі зразком. 54

polymorphic поліморфні, “багатоформні”. 29

polymorphism поліморфізм.
 ad hoc ad hoc. 29
 apparent удаваний. 31
 bounded з обмеженнями. 34
 first-class першого класу. 34, 39
 impredicative непередикативний. 42
 inclusion включний. 31
 inheritance спадковий. 31
 parametric параметричний. 29
 predicative предикативний. 42
 rank- n рангу n . 42
 subtype підтиповий. 31
 universal універсальний. 30

preservation збереження. 19, 20

progress просування. 19

redex редекс. 22, 65, *see reducible expression*

reducible expression скорочуваний вираз, редукований вираз. 22

refine уточнює. 44

representation представлення. 3

semantic style семантичний стиль. 12

subject reduction скорочення/редукція підмета. *see preservation*

subtype підтип. 31

term терм.

stuck застряглий. 16
typable типізовний. 17
well-typed коректно типізований. 17
tree дерево.
 derivation виведення. 14
type тип.
 algebraic алгебраїчний. 53
 function “стрілочний”. 24
 inductive індуктивний. 53
type checking перевірка типів. 4
 dynamic динамічна. 5
 static статична. 5
type erasure стирання типів. 27
type inference виведення типів. 27
type reconstruction відновлення типів. 27
type safety типова безпека. 19
typing derivation виведення типу. 18
typing relation відношення типізації. 17
value значення. 13, *see term*