

A Gentle Introduction to Mercury

The Functional-Logic Programming Language

Zetian Lin

(Very early draft; October 2025)

© 2025 Zetian Lin. Some rights reserved.
<https://sebastian.graphics>

Licensed under Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International with additional terms. To view a copy of the Creative Commons license, please visit <https://creativecommons.org/licenses/by-nc-nd/4.0/>.

Additional terms requested by the author:

All content is shared according to the policies above with all parties EXCEPT any Wikimedia Foundation projects. Under NO CIRCUMSTANCES may the content of this work be reproduced, distributed or transmitted in any form or by any means, including photocopying, recording, or other electronic or mechanical methods for the use of ANY existing and future Wikimedia Foundation project. Permission for any individual or organization to use the content of this work in any form for Wikimedia Foundation projects is hereby explicitly revoked. Any individual or organization that uses the content for any Wikimedia Foundation project will be subject to legal action.

Contents

Contents

1	Introduction	1
1.1	What is Mercury?	1
1.2	The state of things as of now	2
1.3	The tool we will be using	3
2	A First Taste of Mercury	5
2.1	Prerequisite: logic programming	5
2.2	Writing the program	12
3	Programming in Mercury	27

Introduction

1.1 What is Mercury?

Mercury¹ is a *functional logic* programming language, which is a paradigm that combines features common in certain functional languages (e.g. higher-order constructs, strong advanced typing) and logic programming semantics (e.g. non-determinism, built-in search). It attempts to combine the clarity and expressiveness of declarative programming with advanced static analysis and error detection features. According to its creators, it's a “general purpose language intended to support the creation of large, reliable and efficient applications.” Although a short-lived renewed interest in this paradigm has been generated by Epic Games's programming language Verse² a few years earlier, it has now seemingly died down unfortunately; the paradigm remains interesting regardless, for the obvious reason that it provides extra edge against both purely functional and purely logic languages.

As are most of the languages that are of similar categories, Mercury has been mostly used for academic/research purposes, but a few commercial products have chosen (or once chosen) Mercury as their main language as well, like PrinceXML³, which is an HTML rendering tool used in publishing.

¹<https://mercurylang.org>

²<https://dev.epicgames.com/documentation/en-us/fortnite/verse-language-reference>

³<https://www.princexml.com>

Other similar languages

A few languages have been said to be similar to Mercury on its different aspects:

- Curry⁴: When discussing functional logic programming, people naturally lean towards a conceptual Haskell-Prolog continuum; if one were to say Mercury is closer to the Prolog side, Curry is clearly closer to the Haskell side of things.
- Clean⁵: Clean is the lesser-known sibling of Haskell famed for its uniqueness typing. It was said that the ability to explicitly specify uniqueness is not unlike that of Mercury.

1.2 The state of things as of now

Mercury has been considered a difficult language to learn, even for people who has previous functional language experiences. While Mercury do have unfamiliar aspects from the viewpoint of programmers who have only written in “conventional” languages like Scala or Haskell, I personally find those aspects aren’t as difficult as one may think; I believe, that to a certain point this “difficult” feel comes from other factors, like the lack of a visible software ecosystem:

- Lack of example, tutorials and other introductory text: To be honest, it’s not that we don’t *have* those, it’s simply that most of them tend to be slightly too confusing for common programmers who are used to conventional imperative/functional languages. As of the time of writing (October 2025), we have the following attempts:
 - The tutorial written by Ralph Becket which is listed on Mercury’s official website⁶. Unfortunately this tutorial was left unfinished years ago and its author was said to have since stopped doing Mercury-related things.
 - The tutorials on Mercury’s GitHub repository wiki⁷, which seemingly has ceased major activities since 2018.

⁴<https://curry-lang.org>

⁵<https://clean-lang.org>

⁶<https://mercurylang.org/documentation/learning.html>

⁷<https://github.com/Mercury-Language/mercury/wiki/Tutorial>

- “Learn Mercury By Y Minute”⁸, which didn’t do a very good job at explaining things to a depth a learner would need their materials to have.
 - “Mercury Crash Course” written by J Fondren, on the now-defunct mercury-in.space website⁹, which honestly was too confusing for newcomers to provide any help.
 - A book titled “Declarative Programming in Mercury” is also available by compiling from the TeX source yourself using the source code¹⁰, but this one really isn’t an introductory text...
- Lack of package manager: Same as above, it’s not like we *don’t* have one - in fact we have two, one called Merchant by Stewart Slocum¹¹, and one called **mmc-get**¹² by J Fondren, which used to have a handful of packages¹³. They are not compatible with each other (of course), which I doubt have ever created any kind of problem since nearly no one used them.
 - Lack of libraries: This is the real pain point. If a language did not have enough libraries, people are much less likely to use the language to make things; and if not enough people are using the language, there wouldn’t be enough demand and thus there wouldn’t be enough libraries. It’s a vicious cycle that requires a *lot* of effort to break.

But all of these situations won’t change if we simply do nothing and hope for that one dumbass to put their love and effort for little to none rewards, which is why I have decided to become that dumbass myself and write this book to bridge the gap somewhat.

1.3 The tool we will be using

As of now (October 2025) there is only one implementation of Mercury: the Melbourne Mercury Compiler. All code in this book has been tested against

⁸<https://learnxinyminutes.com/mercury/>

⁹You can still read it on Wayback Machine: <https://web.archive.org/web/20220815044519/https://mercury-in.space/crash.html>

¹⁰Available on GitHub: https://github.com/Mercury-Language/books/tree/master/decl_prog

¹¹<https://github.com/stewy33/merchant>

¹²<https://github.com/jrfondren/mmc-get>

¹³<https://web.archive.org/web/20221130042103/https://mercury-in.space/packages/packages.list>

version 22.01.8. Installing the compiler is probably one of the most painful part: since there is no official prebuilt binaries, you will have to download the source code, compile it and install it yourself. If you're using an operating system with a package manager (like most Linuxes), their package repositories might have packages that ease the complexity of managing a little bit, but you would still need to compile the compiler, which takes a very long time.

A First Taste of Mercury

Unfortunately, due to Mercury currently not having a read-eval-print-loop, it's nigh impossible to properly introduce the basic part of Mercury in details without introducing some kind of scaffolding first. For this reason, in this chapter we'll be forcing our way into writing and completing a tiny utility program; after this, we will have a basic idea of how to write a bare-bone command line program in Mercury, which we can then use as a starting point for exploring the language.

2.1 Prerequisite: logic programming

(You can skip this part if you have learned Prolog before – in fact, I might even go as far as to tell you to learn Prolog instead of reading this chapter, which is my poor attempt at stuffing as much essential knowledge within as short a span as possible. You don't have to learn a lot – basic Prolog, like the amount taught by *Learn Prolog Now!*¹, is more than sufficient.)

While it might not seem obvious from an outsider's view, logic programming did in fact come from first-order predicate logic; more specifically, it comes from a specific interpretation of a specific class of first-order predicate formula. This class of formula, called **Horn Clause** after the logician Alfred Horn, can be one of the following forms:

¹<https://lpn.swi-prolog.org/lpnpage.php?pageid=online>

- “Rule”, which takes the form of $Q \leftarrow P_0 \wedge P_1 \wedge \dots \wedge P_n$. We will read this as “ Q holds when all of P_0 and P_1 and the rest – all the way to n – holds”
- “Fact”, which takes the form of $P \leftarrow \top$. We will read this as “ P holds (regardless)”. Sometimes it’s also shortened further into P . in some literature.
- “Goal”, which takes the form of $\perp \leftarrow P_0 \wedge P_1 \wedge \dots \wedge P_n$. We will read this as “Show all of P_0 and P_1 and the rest holds”.

Within these formula, all the Q and P_n have the form of $f(arg_0, arg_1, \dots arg_n)$ where f is simply a name and all the arg_n are themselves first-order predicate logic terms. If how this has anything to do with writing programs isn’t clear to you, it probably will after reading the definition of the factorial function, written in Horn clauses:

$$\begin{aligned} fac(0, 1) &\leftarrow \top \\ fac(1, 1) &\leftarrow \top \\ fac(N, M) &\leftarrow minus(N, 1, N_0) \wedge fac(N_0, M_0) \wedge mult(M_0, N, M) \end{aligned}$$

Compare this with the following:

$$\begin{aligned} fac(0) &= 1 \\ fac(1) &= 1 \\ fac(N) &= M \textbf{ where } N_0 = N - 1, M_0 = fac(N_0), M = N \times M_0 \end{aligned}$$

...and the following code in Python:

```

1 def fac(n):
2     if n == 0: return 0
3     elif n == 1: return 1
4     else:
5         n0 = n - 1
6         m0 = fac(n0)
7         return n * m0

```

The similarity should be obvious at this point; in a very inaccurate sense, you can think of it as:

- Rules as function definition;
- Predicates as function calls;
- Conjunction as sequencing;

- Facts as base cases;
- Goals as function calls & entry points;

This, of course, is not all there is to it. We've covered the representation, we shall now cover the reduction.

Unification

One thing to keep in mind is that theoretically there is no “computation” (as in calculation) involved; what we have (mostly) is a thing we call *unification*.² Unification roughly means to check if two terms “match”.

All terms are of one of the following forms:

- Variables, which we will denote with words that starts with an uppercase letter;
- Atoms, which we will denote with words that starts with a lowercase letter;
 - (This name, despite being the commonly used name for languages like LISP or Erlang, might come as obscure for other people; a more descriptive name would probably be “scalar values”.)
- Functors³, which is of the form $f(a_0, \dots, a_n)$, in which f is an atom and all the a are terms (which can themselves be variables or atoms or functors, etc.)

The unification of two terms thus can be said as follows:

- A variable matches with anything (unless it forms recursive reference; see below);
- Atoms only match with themselves, e.g. the atom *apple* would only match with *apple* and will not match with *pear*. (But obviously it would still match with variables, since variables match with anything)

²To some people this (and resolution, which we will discuss later) looks like a rewriting system, and some people do claim that Prolog is just a rewriting system with backtracking. As someone who've worked with rewriting systems before, I'm not going to comment on their opinions, but they do be looking similar.

³The name “functor” has been used to refer to a handful of different concepts; in the context of logic programming, especially languages that's derived from Prolog, you can simply think of it as “some data that's shaped like a function call”.

- Functors match if and only if all the following conditions are met:
 - The “head” (or “the functor’s name”) matches, e.g. $f(A, B)$ matches with $f(a, b)$ but not $g(a, b)$ because the latter is g instead of f .
 - The arity (i.e. the number of arguments) matches, e.g. $f(A, B)$ matches $f(a, b)$ but not $f(c, d, e)$ despite both being f ; the former has 2 arguments while the latter has 3. The convention is to denote the arity after the functor’s name, e.g. $f/2$ and $f/3$.
 - All of the arguments match with each other, e.g. $f(a, B)$ matches with $f(a, b)$ and $f(a, c)$ but not with $f(b, a)$ despite both being $f/2$; the first arguments, namely a and b , does not match.

The action of unification results in a mapping from variables to terms; (one could say that this is how “value assigning” happens, although that really isn’t the accurate wording.) A common name for this mapping would be “subst”, short for substitution. The observant might want to ask if one can match a variable with a term that contains the same variable either directly or indirectly (e.g. A with $f(A)$). Theoretically this is not valid because one can never have a fully grounded solution for such a subst, but whether an occurrence of such case will trigger some kind of runtime panic depends on the language and its implementation. Some languages (like Prolog) permits it (even if this kind of situations almost certainly end in a failure state), some languages (like miniKanren) performs strict occur check and does not allow any such occurrences.

Some of you may ask if pattern matching can also be seen as a form of unification. I’ve met people online who were very pedantic about the name and fervently insisted that pattern matching is not unification; almost all of those kind of people are doing this because they want to be perceived as more sophisticated than they truly are so they can look down upon other people more safely. I would say one can think of it as somewhat similar; one must understand that pattern matching (and the subsequent binding of variables) is a “one way thing”, e.g. unifying $f(a, J)$ with $f(K, b)$ would result in a subst of $J = b, K = a$, but with pattern matching one side is always a fully grounded term, if we can even consider values as terms in this sense.

Resolution

In this section I would only introduce one way of resolution. It is – I believe – conceptually the simplest kind of resolution. There probably are other

kind of resolution technique out there, but they are out of scope for this tutorial.

Remember that a goal is of the form $\perp \leftarrow P_0 \wedge \dots \wedge P_n$ (we'll omit the \perp part for brevity from now on) and act as entry points in our interpretation; naturally, if there's no more P , the goal is empty, and the computation halts.

The simplest method of resolution is thus as follows. For the goal $\leftarrow P_0 \wedge \dots \wedge P_n$, we take the left-most subterm, in this case P_0 ; we attempt to find among all the previously defined rules and facts the one that can successfully unify with P_0 (instantiated according to the current subst). If we couldn't find one, then the current attempt is considered a failure (what happens after the failure depends; see later section on cut); If we were successful in finding one, then two things happen:

- Because we performed a unification, we now have a subst;
- And because we picked a rule/fact, we now have a “right-hand side”, which could be empty or containing other goals.

This subst resulted from the unification is combined with all the subst we've had so far, and the “right-hand side” of the rule/fact we found would be appended to the left-hand side of the current goal.

Let's look at an example. Assume we have the following program in Horn clauses:

```

parent(tom, john) ← ⊤
parent(tom, evans) ← ⊤
parent(tom, sarah) ← ⊤
sibling(A, B) ← parent(C, A) ∧ parent(C, B) ∧ not(equal(A, B))

```

...and the goal we're asking is $\leftarrow sibling(john, A)$; the following things happen:

- $sibling(john, A)$ matches with $sibling(A, B)$, (but we cannot directly say that $john = A, A = B$, because the two A should clearly mean different things! with a bit of renaming, we shall say the goal is $sibling(john, A_0)$).
- The unification of the step above results in the subst $john = A, A_0 = B$.

- It contains a right-hand side of positive length; we shall add them to the goal with respect of the subst. The goal is now $\leftarrow \text{parent}(C, \text{john}) \wedge \text{parent}(C, A_0) \wedge \text{not}(\text{equal}(\text{john}, A_0))$.
- We now try to find a rule/fact that can unify $\text{parent}(C, \text{john})$, and we found $\text{parent}(\text{tom}, \text{john})$. This gives us $C = \text{tom}$. Since it's a fact, there is no right-hand side to be added. The goal is now $\leftarrow \text{parent}(\text{tom}, A_0) \wedge \text{not}(\text{equal}(\text{john}, A_0))$.
- We now try to find a rule/fact that can unify $\text{parent}(\text{tom}, A_0)$, and we again found $\text{parent}(\text{tom}, \text{john})$. This gives us $A_0 = \text{john}$.
- The goal, which is now $\leftarrow \text{not}(\text{equal}(\text{john}, \text{john}))$ with respect to the current subst, fails. If we have committed to this choice (by cut or by other means), the entire thing fails; but we didn't, thus we backtrack. The goal is now once again $\leftarrow \text{parent}(\text{tom}, A_0) \wedge \text{not}(\text{equal}(\text{john}, A_0))$, and we attempt to find the next thing that can unify with $\text{parent}(\text{tom}, A_0)$.
- We found $\text{parent}(\text{tom}, \text{evans})$, which results in the subst $A_0 = \text{evans}$.
- The goal is now $\leftarrow \text{not}(\text{equal}(\text{john}, \text{evans}))$, which would succeed without a subst (or with an empty subst, depending on how you see it), leaving the goal empty. Since the goal is now empty, we consider whatever we have now is a valid solution. The subst for this solution would be $A = \text{john}, C = \text{tom}, A_0 = \text{evans}$, and the solution itself would be $\text{sibling}(\text{john}, \text{evans})$.
- If more result is requested, the latest unification and its subst would be discarded and a new instance of $\text{parent}(\text{tom}, A_0)$ would be looked for, which would result in the subst $A_0 = \text{sarah}$ and the solution $\text{sibling}(\text{john}, \text{sarah})$.
- It's easy to see that we can no longer find any new solutions after this point, so the resolution for the goal ends here.

Cut

Cut is mainly a Prolog thing; when triggered, it forces the resolution to commit to all the choices it has made up till that point, artificially cutting the backtracking tree. This is one of the major reasons why Prolog is being considered as “not so pure”. Similar constructs exists in other logic

programming languages (e.g. Mercury's `cc_multi` and `cc_nondet`), albeit they may work in a fundamentally different manner.

Using our program above as an example; if we write the first fact as $\text{parent}(\text{tom}, \text{john}) \leftarrow !$ instead of $\text{parent}(\text{tom}, \text{john}) \leftarrow \top$ (the exclamation mark is the "cut operator" in Prolog), we wouldn't have any solutions, since by the time we reached the goal $\leftarrow \text{not}(\text{equal}(\text{john}, \text{john}))$, we have already committed to the choices of *john* and discarded other possibilities.

Disjunction

Other than conjunction which can be interpreted as sequencing, one can also use disjunction in Horn clauses, which can be interpreted as branching. The resolution of a disjunct goal naturally depends on the resolution of its subgoals, and any solution is a solution of the overall disjunct goal if it is also a solution of any of its subgoals. The resolution of a disjunct goal starts from one of its subgoals; if that subgoal fails, it moves onto the next subgoal, until any one of these subgoals succeed (which results in an overall success) or all of them failed (which results in an overall failure). The order in which each branch gets resolved theoretically shouldn't matter, but in reality different order would lead to different backtracking tree, which in languages with cut (e.g. Prolog) and similar constructs could lead to order-dependent results.

Unification, part 2

In short, due to how unification works, you can do certain thing "in reverse". For example, assume that we have a predicate `append/3` which, at surface, is something you append two lists together:

$$\begin{aligned} \text{append}(\text{nil}, \text{nil}, \text{nil}) &\leftarrow \top \\ \text{append}(\text{nil}, A, A) &\leftarrow \top \\ \text{append}(\text{cons}(H, T), A, \text{cons}(H, R)) &\leftarrow \text{append}(T, A, R) \end{aligned}$$

If you query goals where the first two arguments are "solid" terms, you will get the result bound to the variable at the third argument; but at the same time, you can:

- Call the predicate with three "solid" lists, with which you can check if appending the first two lists does end up with the same value as the third one. If that *is* the case, then you'll eventually hit one of the

first two facts and have an overall success; if that *isn't* the case, you'll have a failure.

- Call the predicate with two variables and one “solid” list, with which you have a goal that asks “which two lists can be appended and end up with the same value as this third list”.

For the second one, we'll discuss the case *append(A, B, cons(c, nil))* here; you can try to resolve a few examples on your own as an exercise. In the case of *append(A, B, cons(c, nil))*, the resolution would first match with the second fact (which will give $A = nil, B = cons(c, nil)$) and then it would match with the third rule (which, if you work it out, will give $A = cons(c, nil), B = nil$.)

2.2 Writing the program

In this part of the chapter, we will try to write a program which I called **rall**, which stands for “adding carriage return to **all** lines”. From this “full name” you can tell it is a very small utility that converts Unix-style newlines (LF) to Windows-style newlines (CRLF). It covers some of the most basic things when it comes to software development, and I reckon it would serve as a nice beginner (in terms of language, not programming in general) project.

File name / module name

Each Mercury source file would be a module, whose name should be the same as the file name (without the extension name part). Each module would need to have at least one **interface** section and one **implementation** section.

For our program, we would save the file under the name **rall.m**, thus the module name would be **rall**.

```

1  :- module rall.
2  :- interface.
3  :- implementation.
```

Entry point

With Mercury, the entry point of the program (e.g. **main** in C) must be defined as a **main/2**, i.e. **main** with 2 parameters. If you attempt to define **main** as anything else, you would get a compile error:

```

1 % mmc ./rall.m
2 /usr/bin/ld: rall_init.o: in function `mercury_init':
3 rall_init.c:(.text+0x4a0): undefined reference to `<predicate '
  main'/2 mode 0>'
4 collect2: error: ld returned 1 exit status

```

The two parameters of this `main/2` must be of type `io.state`. If you don't define it as `io.state`, you would get a compile error:

```

1 ./rall.m:005: Error: `main'/2 must have mode `(di,uo)'.
2 ./rall.m:005: Error: both arguments of `main/2' must have type `io
  .state'.

```

To be able to use `io.state`, we should first import the `io` module, which is done with the `import_module` declaration.

So we now have the following:

```

1 :- interface.
2 % NOTE: to use io.state we must import the io module.
3 :- import_module io.
4 :- pred main(io.state, io.state).

```

You might wonder what exactly is `io.state` and why `main` needs two of them; we will explain this later. At the very least, these are **not** the same as `argc` and `argv` which you would see in C. There are dedicated predicates for retrieving command line arguments; we will get to that later.

The `(di,uo)` part needs explanation. The message says it's a **mode**. What is a mode? Modes are things that describe the change of instantiatedness before and after the execution of a goal. Instantiatedness – that's a long word. What is instantiatedness? Roughly speaking, it's a concept for describing whether a slot is free or bound. (In Mercury, instantiatedness is actually a tree, because obviously you'd have – or at the very least you could manually construct – terms that are only bounded at certain parts.) Two kinds of basic instantiatedness exists in Mercury: **free**, which refers to variables that are not bound by any values; **bound**, which means that the term is not a variable but rather some concrete term at least at that level. Consider the term `blah(A, b)`; it's **bound** at the top (with `blah`) and at one of the children node (at `b`) but **free** at another children node (at `A`, provided that `A` does not have a value already). A mode, **ground**, also exists, and refers to the terms that are *completely* bound (e.g. `blah(a, b)` is **ground** but `blah(A, b)` isn't, provided that `A` is still **free**.)

With this knowledge, we should first understand what **in** and **out** modes actually are. If you have worked with some other languages you might have seen thing similar to this before: C libraries often ask you to pass in a reference for retrieving the actual result because the return value is used for returning error status; you can think of it as “in” (inputting) parameters

and “out” (returning) parameters; some languages (e.g. Ada) even explicitly label them as such. I do not wish to introduce you to wrong analogies that will become detrimental for your future learning, but I have to say they do be somewhat similar.

That said, **in** and **out** is properly defined in Mercury, as follows:

```
1 :- mode in == ground >> ground.
2 :- mode out == free >> ground.
```

This largely fits with our intuition about in and out parameters: we expect **in** arguments to be fully grounded and they stay grounded, and we expect **out** arguments are variables which we will provide a solid term to (thus making it grounded).

There are also “polymorphic” version of **in** and **out**, defined as follows:

```
1 :- mode in(Inst) == Inst >> Inst.
2 :- mode out(Inst) == free >> Inst.
```

You can indeed define your own instantiatedness, with which then you can use these definitions to get the **in** and **out** version of it; but we will not need this for our example this time.

Now we can finally come back to **di** and **uo**. **di** stands for “destructive input”, and **uo** stands for “unique output”. If you have used Rust (or, in the case where you are really adventurous, Clean) before, you might have a vague idea of what this is. In Mercury, there are two special instantiatedness named **unique** and **dead**, the former conceptually refers to values that can only have one reference at all time, and the latter refers to reference that are once “unique” but is now “dead”.

```
1 % unique output - used to create a "unique" value
2 :- mode uo == free >> unique.
3
4 % unique input - used to inspect a unique value without causing
5 % reference to become dead
6 :- mode ui == unique >> unique.
7
8 % destructive input - used to deallocate or reuse the memory
9 % occupied by a value that will not be used.
10 :- mode di == unique >> dead.
```

Up to now, our code would be something that’s similar to this:

```
1 :- module mercury_rall.
2
3 :- interface.
4
5 :- import_module io.
6 :- pred main(io.state, io.state).
7
8 :- implementation.
9
```

```

10  main(_, _) :-
11      % some dummy body for our main predicate.
12      1 = 1.

```

If you compile this, the Mercury compiler would complain about not having a mode declaration. For this reason, we will add the following line and compile:

```

1  :- mode main(di, uo).

```

But this time we would see the compiler complain about different things:

```

1  ./rall.m:007: Error: no determinism declaration for exported
    predicate
2  ./rall.m:007:   `main'/2.
3  ./rall.m:012: In clause for `main(di, uo)':
4  ./rall.m:012:   mode error: argument 2 did not get sufficiently
5  ./rall.m:012:   instantiated.
6  ./rall.m:012:   Final instantiatedness of `HeadVar__2' was `free',
7  ./rall.m:012:   expected final instantiatedness was `unique'.

```

The compiler expects the second argument would be supplied with a **free** argument and that argument should become a **unique** value at the end of **main**! How do we do such a thing? Fortunately the **io** module has what we need to stuff the body of **main**, and we will use one of them:

```

1  main(In, Out) :-
2      io.write_string("blah", In, Out).

```

But this time the compiler started complaining about other things:

```

1  ./rall.m:007: Error: no determinism declaration for exported
    predicate
2  ./rall.m:007:   `main'/2.

```

Determinism in this case, informally speaking, refers to “how a certain thing would succeed/fail”. It’s not something we’d care about in other languages, at least not in an explicit, supported-by-the-language-itself manner; (we normally only talk about when programs terminate at a certain state or not.) In Mercury, we have the following determinism categories:

- Deterministic (**det**): guaranteed to have one and exactly one solution.
- Semideterministic (**semidet**): have exactly one solution, but does not guarantee to produce it.
- Multisolution (**multi**): guaranteed to have a solution among possibly many solutions.
- Nondeterministic (**nondet**): have possibly many solutions, does not guarantee to produce one.

- Failure (**failure**): cases where there's zero solutions. They are not actual errors but a part of the logic (e.g. arity mismatch during unification, which will never produce a solution because the arity is different).
- Errorneous (**errorneous**): also have zero solutions, but they **do** represent actual errors which in other languages would be represented in the form of runtime exception throw or panic.

Some people might not understand why **det** and **semidet** are separate things. Imagine a function that takes the "head" of a linked list; this function is obviously only defined on non-empty list and not defined on empty lists, so for any given list there's either only one solution or no solution, thus **semidet** instead of **det**.

The difference between **errorneous** and **failure** might be clearer if I explain their behaviours when it comes to negation. Basically in Mercury (and other logic programming languages), you can take a predicate, put a negation on it, and then ask for a case where it **does not succeed**. Naturally, the negation of determinisms that are guaranteed to produce a result like **det** and **multi** would be **failure** (i.e. negating a definite success would be a failure), the negation of **failure** would be **det** (i.e. negating a definite failure would be a definite success), and the negation of **semidet** and **nondet** would be **semidet** (i.e. negating these cases turns it into a case that asks if the match would be successful or not); but negating an **errorneous** would only produce an **errorneous**.

(If you're going to ask this question – yes, solving the problem of perfect determinism inference does mean solving the halting problem, and is thus undecidable.)

Other than these six categories there are also this thing called the "committed choice nondeterminism", which adds two more determinism class: **cc_multi** and **cc_nondet**. The difference between committed choice nondeterminism and "normal" nondeterminism is that normal nondeterminism supports backtracking while CC nondeterminism, despite potentially having more than one solutions, commits to only one of them and thus does not backtrack. The **main** entry point can be defined to be **det** or **cc_multi**, since both of them guarantee one and only one solution (**det**, of course, is more strict than **cc_multi**, and if the Mercury compiler can determine something that should be able to be a **det** got labelled as a **cc_multi**, it would spit out a warning saying you could've gone with the stricter option.)

In this case, we should add the string **is det** at the end of our mode declaration, so that the whole line would be **:- mode main(di, uo) is**

det.. After this modification, the compiler should have finally stopped complaining and gives you an executable; when you run it, it would display a string **blah**, which should be obvious. At this point, the code shall look like something like this:

```

1  :- module rall.
2
3  :- interface.
4
5  :- import_module io.
6  :- pred main(io.state, io.state).
7  :- mode main(di, uo) is det.
8
9  % also: you can combine the `pred` and the `mode` line into one
10 % like this:
11 %      :- pred main(io.state::di, io.state::uo) is det.
12
13 :- implementation.
14
15 main(In, Out) :-
16     io.write_string("blah", In, Out).
```

Now, if you attempt to write multiple strings like this:

```

1  main(In, Out) :-
2      io.write_string("blah", In, Out),
3      io.write_string("blah", In, Out).
```

The compiler would produce this error:

```

1  ./rall.m:013: In clause for `main(di, uo)':
2  ./rall.m:013:   in argument 2 of call to predicate `io.
   write_string'/3:
3  ./rall.m:013:   unique-mode error: the called procedure would
   clobber its
4  ./rall.m:013:   argument, but variable `In' is still live.
5  For more information, recompile with `-E'.
```

If we look up the declaration of **io.write_string**:

```

1  :- pred write_string(string::in, io::di, io::uo) is det.
```

You should already know that **di** specifies a turn of a **unique** value into a **dead** value, so **In** after the first **write_string** would be considered **dead** and not **unique**, which does not fit the requirement of the **di** of the second **write_string**. If you try to experiment and do this:

```

1  main(In, Out) :-
2      io.write_string("blah", In, Out),
3      io.write_string("blah", Out, Out2).
```

This would have the following violation: since we declared **Out** to be **uo** which is **free >> unique**, **Out** must be **unique** at the end of **main**, but the

second `write_string` turns `Out` into a `dead` because that slot is `di` which is `unique >> dead`. A solution to this would be to do this:

```
1  main(In, Out) :-
2      io.write_string("blah", In, Out1),
3      io.write_string("blah", Out1, Out).
```

This would compile, and the generated executable would write "blah" twice, as expected.

To keep coming up with new variable name is tedious if this gets long. To solve this problem, Mercury have something called the "state variable". With state variables, the example above can be written like this:

```
1  main(!.IO, !:IO) :-
2      io.write_string("blah", !.IO, !:IO),
3      io.write_string("blah", !.IO, !:IO).
```

`!.IO` refers to the value bound to `IO` at the current moment, and `!:IO` makes the value at that slot bound to `IO`; so in the first `write_string`, the *current* value `IO` was destroyed, and the value bound to `Out1` in our non-state-variable version was bound to `IO` as its next value; and in the second `write_string`, the value of `Out1`, which itself has become the current value of `IO` and being referred to with `!.IO`, was destroyed, and the value bound to `Out` in our non-state-variable version was bound to `IO` as its next value, and subsequently became the output of `main`.

You can write `!.IO, !:IO` as `!IO`, because the former is also quite tedious as well:

```
1  main(!IO) :-
2      io.write_string("blah", !IO),
3      io.write_string("blah", !IO).
```

It would still feel a bit tedious for people who are used to other programming languages, but I suppose this is what you give up for being as explicit as possible for the sake of improved correctness and things...

Command line arguments

The predicate we need for retrieving the command line arguments is also in the `io` module:

```
1  main(!IO) :-
2      io.command_line_arguments(Argv, !IO),
3      % ...
```

This would bound the actual value to `Argv`. From the definition of `command_line_arguments` we can know that the type of `Argv` would be a `list(string)` – a list of strings. It would be necessary to import the `list`

module from now on. This value does not contain the name of the program (for which `io` module has other predicates to retrieve). It would also be useful to check if the length of `Argv` is exactly 2:

```

1 :- import_module list.
2 % ...
3 main(!IO) :-
4     io.command_line_arguments(Argv, !IO),
5     ( if length(Argv) \= 2
6       then (
7           % you can specify which stream you are writing to btw
8           io.write_string(
9               io.stderr_stream,
10              "usage: rall [inputfilename] [outputfilename]\n",
11              !IO)
12         )
13     else (
14         % ...
15     )
16 ).

```

(Note that if you have previous experience with Prolog, the condition for the `if` construct might be surprising, since in Prolog, `=` unify things *structurally*; for example, the goal `3 + 4 = 7` would fail in Prolog despite 3 plus 4 do equal to 7, since `3 + 4` is a functor (interpreted as ``+(3,4)`) and 7 is an atom; this is different in Mercury, where things *would* be evaluated. Also, `length` in Mercury is defined both as a predicate and as a **function**, and in this case the function got evaluated as it's defined in the standard library resulting an actual number.

The `then` clause is simple – we write the usage string and call it done:

```

1 :- import_module list.
2 % ...
3 main(!IO) :-
4     io.command_line_arguments(Argv, !IO),
5     ( if length(Argv) \= 2
6       then (
7           % you can specify which stream you are writing to btw
8           io.write_string(
9               io.stderr_stream,
10              "usage: rall [inputfilename] [outputfilename]\n",
11              !IO)
12         )
13     else (
14         % ...
15     )
16 ).

```

Now we extract the command line arguments. `list` module contains `index0` and `index1`, which are predicates/functions that uses 0-based index and 1-based index respectively; but these aren't the ones we would choose to use here because they are **semidet**. (Exercise: can you figure out the

reason why they're `semidet` instead of `det` on your own?) We use their `det` counterparts: `det_index0` and `det_index1`, which will throw runtime exceptions when failing:

```

1 :- import_module list.
2 % ...
3 main(!IO) :-
4     io.command_line_arguments(Argv, !IO),
5     ( if length(Argv) \= 2
6       then (
7           io.write_string(
8               io.stderr_stream,
9               "usage: rall [inputfilename] [outputfilename]\n",
10              !IO)
11         )
12       else (
13           InputFilePath = det_index0(Argv, 0),
14           OutputFilePath = det_index0(Argv, 1),
15           rall(InputFilePath, OutputFilePath, !IO)
16           % ...
17       )
18     ).

```

By the way, if we import module simply by doing `:- import_module name`, We don't need to fully specify the module name, as long as there's no ambiguity. Naturally, the good practice for a lot of common cases is to fully specify; but some editors (Emacs, for example) might implement some kind of "smart indent" algorithm which could put the cursor at undesirable positions and be a nuisance. In cases like this, you might find writing code this way to be more comfortable.

It would be also be a good time to write our first predicate that is not the entry point...

Opening & closing files for input/output

In the last code snippet we've write `rall(InputFilePath, OutputFilePath, !IO)`. Despite this looks like `rall/3`, but it's secretly `rall/4`, and we should write our declaration as such:

```

1 :- pred rall(string, string, io.state, io.state).
2 :- mode rall(in, in, di, uo) is det.

```

`io` module has `open_input/4` and `open_output/4`, both of which are `det`:

```

1     rall(InputFilePath, OutputFilePath, !IO) :-
2         io.open_input(InputFilePath, RIn, !IO),
3         % ...

```

The declaration of `open_input` is as follows:

```

1  :- pred open_input(
2      string::in, io.res(io.text_input_stream)::out,
3      io::di, io::uo) is det.

```

We know that whatever variable we put at the second slot would be the input stream we want (or in the case of error, whatever error we could catch). Its value is of type `io.res`, whose definition is as follows:

```

1  :- type res(T)
2  ---> ok(T)
3  ;    error(io.error).

```

This is how Mercury defines algebraic datatypes. If you have previous experience with Haskell, OCaml, Rust or similar languages, this should be natural for you.

Pattern matching is quite simple – you just do `Var = ok(OkValue)` and `Var = error(ErrorValue)` like this:

```

1  % ...
2  ( Var = ok(OkValue),
3  % `ok` clause
4  );
5  ( Var = error(ErrorValue),
6  % `error` clause
7  )

```

Since we combined them with a disjunction, when the unification in the first branch (in this case, `Var = ok(OkValue)`) fails, the program flow would go to the second branch, and if the unification there fails (and there's a third branch), it would go to the third clause, etc.; and if all branches fail, the whole thing fails (and of course we'd expect the Mercury compiler gives us a determinism verdict of something like `nondet`, but I digress). We'll handle all the cases a value of type `res(T)` could have, and if we do that, the compiler should recognize that we'll not fail as long as both of our `ok` clause and `error` clause don't have a determinism verdict that implies a possibility of failing.

The definition of `io.error` is as follows:

```

1  :- type io.error.    % Use error_message to decode it.

```

This type is opaque; we need `error_message` to handle it. Its definition is as follows:

```

1  % Look up the error message corresponding to a particular error
   % code.
2  :- func error_message(io.error) = string.
3  :- pred error_message(io.error::in, string::out) is det.

```

Now we know we can obtain a **string** from an **io.error**; we can decide either to **throw** (which would cause the program to exit prematurely) or to write it to standard error (and try to exit afterwards). In this example we use **throw** because it's simpler:

```

1  % throw/1 requires the exception module.
2  :- import_module exception.
3
4  % ...
5  ( RIn = error(ErrorValue), throw(error_message(ErrorValue)) );
```

Before we move on with it, I must also explain something else. This isn't the case in Prolog, but in Mercury calls are "curried". Using the declaration of **error_message** above as an example, the term **error_message(Var)** *matches both* **func error_message/1** and **pred error_message/2**, and would thus be judged as having the type **string** from **func error_message/1** and the type **pred(string)** from **pred error_message/2** (which is a **pred** that takes a **string** as its argument) *at the same time*, i.e. a type ambiguity. In Prolog it will only match **error_message/2**, but in Mercury this is how it is...

To resolve this you need to explicitly state the type you want by writing **{type}**; in this case, we'll choose the **func** one and specify it as **string**, just like this:

```

1  %                                     here, the ":string" part. ---v
2  ( RIn = error(ErrorValue), throw(error_message(ErrorValue):string)
   );
```

If you don't do this, you'll get a compile error like this:

```

1  ./rall.m:025: In clause for predicate main'/2:
2  ./rall.m:025:   warning: variable 'OutputFilePath' occurs only once
3  ./rall.m:025:   in
4  ./rall.m:031:   this scope.
5  ./rall.m:031: In clause for predicate rall'/3:
6  ./rall.m:031:   error: ambiguous overloading causes type ambiguity
7  .
8  ./rall.m:031:   Possible type assignments include:
9  ./rall.m:031:   V_20:
10 ./rall.m:031:   pred(
11 ./rall.m:031:   string
12 ./rall.m:031:   )'
13 ./rall.m:031:   or
14 ./rall.m:031:   string'
15 ./rall.m:031:   You will need to add an explicit type
16 ./rall.m:031:   qualification to
17 ./rall.m:031:   resolve the type ambiguity. The way to add an
18 ./rall.m:031:   explicit
19 ./rall.m:031:   type qualification is to use "with_type". For
20 ./rall.m:031:   details see
21 ./rall.m:031:   the "Explicit type qualification" sub-section of
22 the
```

```

17 ./rall.m:031:  "Data-terms" section of the "Syntax" chapter of
    the
18 ./rall.m:031:  Mercury language reference manual.

```

Now we can have the following code. Opening an output stream is a near-identical process:

```

1 % throw/1 requires module `exception`.
2 :- import_module exception.
3 % ...
4
5 :- pred rall(string, string, io.state, io.state).
6 :- mode rall(in, in, di, uo) is det.
7 rall(InputFilePath, OutputFilePath, !IO) :-
8     open_input(InputFilePath, RIn, !IO),
9     ( ( RIn = error(ErrorValue), throw(error_message(ErrorValue):
    string) );
10       ( RIn = ok(InputStream),
11         open_output(OutputFilePath, ROut, !IO),
12         ( ( ROut = error(ErrorValue), throw(error_message(
    ErrorValue):string) );
13           ( ROut = ok(OutputStream),
14             process_stream(InputStream, OutputStream, !IO)
15           )
16       ),
17       close_output(OutputStream, !IO),
18       close_input(InputStream, !IO)
19     )
20 ).
21

```

Actually processing the data

We'll do it as follows: each time we read a single character from the input stream, we check if it's a line feed; if it is, we write a carriage return *and* a line feed to the output stream; or else, we simply write that character to the output stream; we repeat this until we've hit the EOF of the input stream. `read_char` produces a different kind of result named **result**, which has a separate constructor **eof**. You should know how to handle this by now...

Character literals are enclosed with single quotes, by the way:

```

1 :- pred process_stream(
2     io.text_input_stream,
3     io.text_output_stream,
4     io.state, io.state).
5 :- mode process_stream(in, in, di, uo) is det.
6 process_stream(InStream, OutStream, !IO) :-
7     read_char(InStream, InRes, !IO),
8     ( ( InRes = eof );
9       ( InRes = error(Error), throw(error_message(Error):string) );
10      ( InRes = ok(Char),
11        ( if Char = '\n'

```

```

12         then (
13             write_char(OutputStream, '\r', !IO),
14             write_char(OutputStream, '\n', !IO)
15         )
16         else (
17             write_char(OutputStream, Char, !IO)
18         )
19     ),
20     process_stream(InStream, OutputStream, !IO)
21 )
22 ).
23

```

Full program listing

The following code is the listing of the completed **rall** program:

```

1 :- module rall.
2 :- interface.
3 :- import_module io.
4
5 :- pred main(io.state, io.state).
6 :- mode main(di, uo) is det.
7
8 :- implementation.
9
10 :- import_module list, exception.
11
12 main(!IO) :-
13     io.command_line_arguments(Argv, !IO),
14     ( if length(Argv) \= 2
15         then (
16             io.write_string(
17                 io.stderr_stream,
18                 "usage: rall [inputfilename] [outputfilename]\n",
19                 !IO)
20         )
21         else (
22             InputFilePath = det_index0(Argv, 0),
23             OutputFilePath = det_index0(Argv, 1),
24             rall(InputFilePath, OutputFilePath, !IO)
25         )
26     ).
27
28 :- pred rall(string, string, io.state, io.state).
29 :- mode rall(in, in, di, uo) is det.
30 rall(InputFilePath, OutputFilePath, !IO) :-
31     open_input(InputFilePath, RIn, !IO),
32     ( ( RIn = error(ErrorValue), throw(error_message(ErrorValue):
33         string) );
34         ( RIn = ok(InputStream),
35             open_output(OutputFilePath, ROut, !IO),
36             ( ( ROut = error(ErrorValue), throw(error_message(
37                 ErrorValue):string) );

```

```

36         ( ROut = ok(OutputStream),
37           process_stream(InputStream, OutputStream, !IO)
38         )
39     ),
40     close_output(OutputStream, !IO),
41     close_input(InputStream, !IO)
42 )
43 ).
44
45 :- pred process_stream(
46     io.text_input_stream,
47     io.text_output_stream,
48     io.state, io.state).
49 :- mode process_stream(in, in, di, uo) is det.
50 process_stream(InStream, OutStream, !IO) :-
51     read_char(InStream, InRes, !IO),
52     ( ( InRes = eof );
53       ( InRes = error(Error), throw(error_message(Error):string) );
54       ( InRes = ok(Char),
55         ( if Char = '\n'
56           then (
57               write_char(OutStream, '\r', !IO),
58               write_char(OutStream, '\n', !IO)
59           )
60           else (
61               write_char(OutStream, Char, !IO)
62           )
63         ),
64       process_stream(InStream, OutStream, !IO)
65     ).
66

```

Listing 2.1: Full code for `rall.m`

It might not be idiomatic Mercury code, but at least it compiles and runs.

We can compile it and produce an executable by calling the command `mmc rall.m`.

Programming in Mercury

Trace Goals

Before we continue we must talk about how to explore things on one's own. A common strategy of debugging in conventional languages would be to insert printing statements at places where one would expect the program to reach with its state containing certain values; this strategy doesn't work well with Mercury, since to perform IO in Mercury you would need to do it through `io.state`, which naturally has the implication of spreading out and forcing everything that uses the predicate to adopt two extra `io.state` parameters.¹

For example, assume we have the following program that prints out the factorial of 5:

```

1 :- module ex01.
2 :- interface.
3 :- import_module io.
4 :- pred main(io::di, io::uo) is det.
5 :- implementation.
6
7 :- import_module int.
8
9 :- pred fac(int, int).
10 :- mode fac(in, out) is det.
11 fac(N, M) :-
12   ( if ( N = 0 ) then ( M = 1 )
13     else if ( N = 1 ) then ( M = 1 )
14       else if ( N < 0 ) then ( M = 1 )
15         else ( fac(N-1, M1), M = N * M1 )

```

¹Which is not unlike the famous “color of function” problem in languages with `async` (e.g. JavaScript)


```

16 )
17 .
18
19 main(!IO) :-
20     io.write_string("Hello, world!\n", !IO),
21     fac(5, M),
22     io.write_int(M, !IO).

```

Let's say we want to add a printing goal in **fac**. Printing, naturally, requires two **io.state** just like we've been mentioning. If we do it by adding **io.state**, we need to perform the following changes:

- Change the **pred** line into `:- pred fac(int, int, io.state, io.state);`
- Change the **mode** line in a similar fashion;
- Change every occurrence of **fac** accordingly;
- Since for every call site of **fac** we now need to pass two extra **io.state** arguments, the environment surrounding the call site should have those extra arguments ready;
 - and for some predicates this could mean adding extra parameters just like what we've done to **fac**;
 - and now the call sites of *those* predicates also need to have their signature changed, and this change would also affect call sites that calls *them*, and the need of change would thus propagate upwards...
- ...and after the debugging is done, we need to change them *back*.

This kind of coding process not only goes against the spirit of Mercury but is also just extremely tedious. Luckily, for this Mercury prepared a feature called **trace goals**.

Trace goals in Mercury are of the format **trace** [*param1*, ...] *Goal*. In such a goal, all variables bound in the parameter list part is available within the *Goal* part. There are many different parameters that one can use, but the most common is probably going to be **io**; by having **io(!IO)** which bound a pair of **io.state** to the state variable of your choice, you can use it to display debug messages in *Goal* freely without having to do the changes listed above. For example, assume we want our program to display the value of each call to **fac** before they return, we can add a trace goal that introduces a pair of **io.state** arguments like this:

```

1 :- module ex01.
2 :- interface.
3 :- import_module io.
4 :- pred main(io::di, io::uo) is det.
5 :- implementation.
6
7 :- import_module int.
8
9 :- pred fac(int, int).
10 :- mode fac(in, out) is det.
11 fac(N, M) :-
12   ( if ( N = 0 ) then ( M = 1 )
13     else if ( N = 1 ) then ( M = 1 )
14     else if ( N < 0 ) then ( M = 1 )
15     else ( fac(N-1, M1), M = N * M1 )
16   ),
17   trace [ io(!IO) ] (
18     io.write_string("fac(", !IO),
19     io.write_int(N, !IO),
20     io.write_string(") = ", !IO),
21     io.write_int(M, !IO),
22     io.write_string("\n", !IO)
23   )
24 .
25
26 main(!IO) :-
27   io.write_string("Hello, world!\n", !IO),
28   fac(5, M),
29   io.write_int(M, !IO).

```

Notice that without changing the signatures of **fac** we can still obtain **io.state** and display things.

Program that uses multiple module

