



HAL
open science

Wasocaml: compiling OCaml to WebAssembly

Léo Andrès, Pierre Chambart, Jean-Christophe Filliâtre

► **To cite this version:**

Léo Andrès, Pierre Chambart, Jean-Christophe Filliâtre. Wasocaml: compiling OCaml to WebAssembly. 2023. hal-04311345v3

HAL Id: hal-04311345

<https://inria.hal.science/hal-04311345v3>

Preprint submitted on 1 Sep 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - ShareAlike 4.0 International License

Wasocaml: compiling OCaml to WebAssembly

LÉO ANDRÈS*, Université Paris-Saclay, CNRS, ENS Paris-Saclay, Inria, Laboratoire Méthodes Formelles, France

PIERRE CHAMBART, OCamlPro SAS, France

JEAN-CHRISTOPHE FILLIÂTRE, Université Paris-Saclay, CNRS, ENS Paris-Saclay, Inria, Laboratoire Méthodes Formelles, France

The limitations of JavaScript as the default language of the web led to the development of Wasm, a secure, efficient and modular language. However, compiling garbage-collected languages to Wasm presents challenges, including the need to compile or re-implement the runtime. Some Wasm extensions such as Wasm-GC are developed by the Wasm working groups to facilitate the compilation of garbage-collected languages. We present Wasocaml, an OCaml to Wasm-GC compiler. It is the first compiler for a real-world functional programming language targeting Wasm-GC. Wasocaml confirms the adequacy of the Wasm-GC proposal for a functional language and had an impact on the design of the proposal. Moreover, the compilation strategies developed within Wasocaml are applicable to other compilers and languages. Indeed, two compilers already used a design similar to ours. Finally, we describe how we plan to handle the C/JavaScript FFIs and effect handlers, in order to allow developers to easily deploy programs mixing C, JavaScript and OCaml code to the web, while maintaining good performances.

CCS Concepts: • **Software and its engineering** → **Compilers**.

Additional Key Words and Phrases: OCaml, Wasm, compilation, value representation, memory, FFI

ACM Reference Format:

Léo Andrès, Pierre Chambart, and Jean-Christophe Filliâtre. 2024. Wasocaml: compiling OCaml to WebAssembly. 1, 1 (September 2024), 29 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 Context

JavaScript is the de facto language of the web. However, it does show its limitations when it comes to performance, security and safety. In order to remedy this, WebAssembly (Wasm for short) [3] has been developed as a secure modular language of predictable performance. Its usage is expanding beyond the web, e.g. finding applications in the cloud (Fastly, Cloudflare) and in the creation of portable binaries.

The first version of Wasm was meant to compile libraries in languages such as C, C++ or Rust to expose them for consumption by a JavaScript program. It was designed with the explicit aim that future versions would cater to the specific needs of other languages and uses. This version of Wasm can be seen as simplified C. Here is an example:

*Also with OCamlPro SAS.

Authors' Contact Information: Léo Andrès, contact@ndrs.fr, Université Paris-Saclay, CNRS, ENS Paris-Saclay, Inria, Laboratoire Méthodes Formelles, Gif-sur-Yvette, France; Pierre Chambart, pierre.chambart@ocamlpro.com, OCamlPro SAS, Paris, France; Jean-Christophe Filliâtre, jean-christophe.filliatre@cirs.fr, Université Paris-Saclay, CNRS, ENS Paris-Saclay, Inria, Laboratoire Méthodes Formelles, Gif-sur-Yvette, France.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

Manuscript submitted to ACM

```

53 (func $fact (param $x i32) (result i32)
54   (if (i32.eq (local.get $x) (i32.const 0))
55     (then (i32.const 1))
56     (else
57       (i32.mul
58         (local.get $x)
59         (call $fact
60           (i32.sub (local.get $x) (i32.const 1)))))))
61
62
63

```

A Wasm module is executed inside what is called an *embedder* or a *host*. In the context of the web, it is a browser and the host language is JavaScript. A Wasm program only manipulates scalar values (integers and floats) but not garbage collected values coming from the host language, such as JavaScript objects. This would be useful in a browser context, but requires some interaction with the GC of the embedder.

Various extensions have been in development for Wasm, some of them being part of the standard already. Our work is based on some of them, that is to say the reference types [20], typed function references [21], GC [18], exception handling [17] and tail calls [19] extensions.

References and GC. The aim of these proposals is to allow the program to manipulate references to different kinds of values such as opaque objects, functions, and garbage-collected values. Since Wasm has to be memory safe, those values are not exposed as raw pointers. Instead, the type system guarantees the proper manipulation of such values.

For instance, in the following example, the `struct` and the `array` can only be manipulated through some primitives such as `struct.get` and `array.get`, there's no way to access their memory address:

```

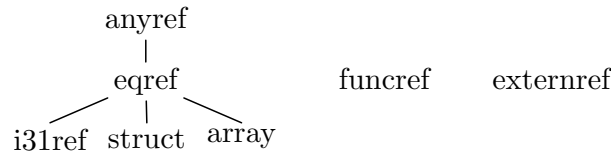
80 (type $f1 (func (param i32) (result i32)))
81 (type $t1 (array i31ref))
82 (type $t2 (struct
83   ((field $a (ref $t1))
84    (field $b (ref $f1))
85    (field $c i32))))
86
87
88
89 (func (param $x (ref $t2)) (result i32)
90   (i31.get_u
91     (array.get $t1
92       (i32.const 0)
93       (struct.get $t2 $a (local.get $x))))))
94
95
96

```

However, Wasm is rather a compilation target than a programming language. It needs to be able to represent the values of any kind of source language. It was not deemed possible to have a type system powerful enough to do that. Instead, the decision was made to have a simple type system, but to rely on dynamic casts to fill the gaps, and guarantee that those casts are efficient.

There is a sub-typing relation that controls which casts are allowed. Some types are predefined by Wasm: `anyref`, `eqref`, `i31ref`, `funcref` and `externref`. Others are user defined: `array` and `struct`. There is a structural sub-typing

relation between user types. Upcasts are implicit, downcasts are explicit and incorrect ones lead to runtime errors. It is possible to dynamically test for type compatibility. The following figure depicts the sub-typing relation:



Contribution to the proposals. The general rule for the Wasm committee is to only include features with a demonstrated use case. As there are currently very few compilers targeting the GC-proposal, some features were lacking conclusive evidence of their usefulness. An example is the `i31ref` type that is not required by the Dart compiler (the only one targeting the GC-proposal at the time). Wasocaml demonstrates the usefulness of `i31ref`. It also validates the GC-proposal on a functional language. We presented Wasocaml to the Wasm-GC working group [2]. It helped in convincing the working group to keep `i31ref` in the proposal.

1.0.1 Exceptions. Another serious limitation of Wasm1 is the absence of an exception mechanism. This is not a problem for compiling either C or Rust, but C++ does require them. It is of course possible to encode exception, but this has a significant runtime cost and limits the interactions with JavaScript (that has exceptions).

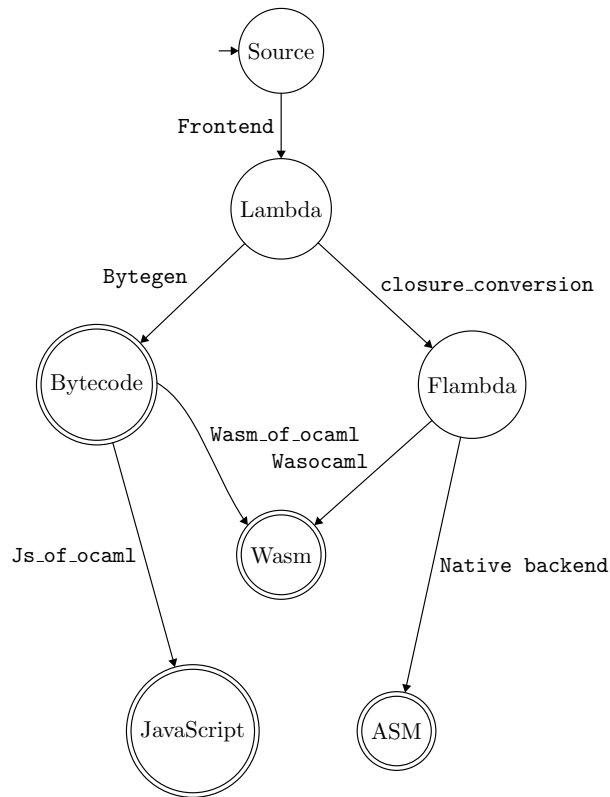
1.0.2 Tail Calls. Wasm1 does not allow tail calls to be optimized. This is a show-stopper when compiling functional languages. In the tail-call proposal, new instructions such as `returncall` guarantee that a tail-call will be optimized.

2 Source Language

What we are interested in is compiling OCaml to Wasm. Before being able to compile, we have to choose where to start from in the OCaml compiler. In the following subsection we will describe the possible choices and then we will give a proper small-step semantics of the chosen language.

2.1 OCaml Intermediate Representation Choice

OCaml has many intermediate representations (IR for short). In the *Lambda* IR, closures are still implicit and the code has not been optimized. The *Bytecode* representation is not well optimized. The *Clambda* one is too low-level for Wasm because code pointers and values are mixed inside closures. *Cmm* is even more low-level: it has pointer arithmetic. The *Flambda* IR is simpler than *Flambda2*, thus it is the one we chose for our first prototype. We will use *Flambda2* in the future, as it provides better optimisations. The following figure depicts a simplified version of the OCaml compilation chain:



2.2 Native OCaml Value Representation

The various OCaml IRs are working on a uniform representation of values. There are two kinds of values: small scalars and heap-allocated blocks.

At runtime, it is possible to test if a value is a scalar or a block. In the native backend, this is done through a technique known as *pointer tagging*, but in the Flambda IR, the precise representation is not imposed yet.

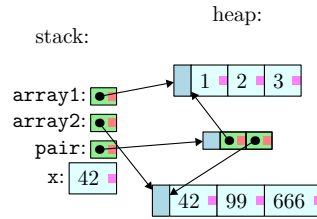
Small scalars are used to represent booleans, unboxed integers, characters or constant constructors of a sum type. Heap-allocated blocks are used to represent arrays, non-constant constructors of a sum type or pairs. For instance, consider the following OCaml code:

```

let array1 = [| 1; 2; 3 |]
let array2 = [| 42; 99; 666 |]
let pair = (array1, array2)
let x = 42

```

In memory, it would be represented as follows:



Note that a given type can have values of both kinds. For instance the empty list is a scalar value but a list with at least one element will be a heap-allocated block.

2.3 Flambda semantics

We give a formal semantics for a subset of Flambda. Compared to the full Flambda it has the following limitations: the objects related primitives have been removed, functions only have one argument, exceptions have been removed but we still have static exceptions that are described later.

2.3.1 Abstract syntax.

Value. A value v is either an integer n or an address a .

$$v ::= n \quad (\text{integer})$$

$$| a \quad (\text{address})$$

An integer represents a value of many types in the OCaml source language such as an `int`, a `bool`, a `char` or the constant constructors of a sum type. An address points to a heap-allocated chunk.

Store. A store \mathcal{M} is a map from addresses to heap-allocated chunks. It is used to represent the state of the memory. Heap-allocated chunks can either be a block made of a tag n and a list of values v^* , a closure $\mathcal{S}.x$, or a set of closures \mathcal{S} . A closure is made of a set of closures \mathcal{S} and of a name x which is its identifier inside the set.

$$\mathcal{M} ::= a \mapsto n, v^* \quad (\text{block})$$

$$| \mathcal{S}.x \quad (\text{closure})$$

$$| \mathcal{S} \quad (\text{set of closure})$$

Set of closures. A set of closures \mathcal{S} represents a set of mutually recursive functions.

$$\mathcal{S} ::= \mathcal{V}, \mathcal{C}$$

It is made of two maps \mathcal{V} and \mathcal{C} . The map \mathcal{V} is the environment (the values captured by the closures). An environment is a map from names to values.

$$\mathcal{V} ::= x \mapsto v$$

The map \mathcal{C} is the closure map (the code of each function). It is a map from names to codes.

$$\mathcal{C} ::= x \mapsto \lambda x x . t$$

261 The code is a function that takes two arguments. The first argument is the real argument of the function. The second
 262 one is the address of the closure being called. It is used to access the environment or to call other functions from the
 263 same set of mutually recursive functions. The body t of the function is a term (see below).
 264

265 *Binary operator.*

267 $op ::= eq$ (equality)
 268 | add (addition)
 270 | sub (subtraction)

274 *Branches.*

275 $\mathcal{B} ::= n \mapsto t$

278 *Term.*

279 $t ::= v$ (value)
 281 | x (identifier)
 282 | $\text{let } x = t \text{ in } t$ (let-binding)
 283 | $x x$ (application)
 284 | $x \leftarrow x$ (mutation)
 285 | $\text{if } x \text{ then } t \text{ else } t$ (conditional)
 286 | $\text{switch } x \ \mathcal{B} \ \mathcal{B}$ (switch)
 287 | $\text{sraise } x \ x^*$ (static raise)
 288 | $\text{scatch } t \text{ with } x \ x^* \ t$ (static catch)
 289 | $\text{while}_t \ t \ \text{do } t$ (while loop)
 290 | $\text{get_field } n \ x$ (block field read)
 291 | $\text{set_field } n \ x \ x$ (block field write)
 292 | $\text{make_block } n \ x^*$ (block creation)
 293 | $\text{project_closure } x \ x$ (closure projection)
 294 | $\text{project_var } x \ x$ (closure environment)
 295 | $\text{move_within_closure } x \ x$ (closure movement)
 296 | $\text{make_set_of_closures } \mathcal{S}$ (closures allocation)
 297 | $\text{pop } t$ (pop environment)
 298 | $x \ op \ x$ (binary operator)

309 Note that the `pop` instruction is an *administrative instruction* that is never present in the concrete syntax coming
 310 from user code.
 311

2.3.2 Small-step semantics.

Configuration. Our semantics is operating on a *configuration*. A configuration is of the form $\mathcal{V}^*, \mathcal{M}, t$ where \mathcal{V}^* is a stack of environments representing the call stack ; \mathcal{M} is a store representing the global memory of the program ; t is the current term being evaluated.

Reduction context.

$$\begin{aligned}
 E ::= & \square \\
 & | \text{let } x = E \text{ in } t \\
 & | \text{while}_t E \text{ do } t \\
 & | \text{scatch } E \text{ with } x \ x^* \ t \\
 & | \text{pop } E
 \end{aligned}$$

Reduction inside a context.

$$\frac{\mathcal{V}^*, \mathcal{M}, t_1 \rightarrow \mathcal{V}'^*, \mathcal{M}', t_2}{\mathcal{V}^*, \mathcal{M}, E(t_1) \rightarrow \mathcal{V}'^*, \mathcal{M}', E(t_2)}$$

Head reductions. In order to reduce cluttering, \mathcal{V}^* and \mathcal{M} are removed from the configuration in the left-hand side of the reduction rules when they are not read and not modified but also from the right-hand side when they are not modified.

$$\text{identifier} \frac{\mathcal{V}_0(x) = v}{\mathcal{V}_0 \mathcal{V}^*, x \rightarrow v}$$

An identifier simply reduces to its image in the current environment.

$$\text{let-binding} \frac{}{\mathcal{V}_0 \mathcal{V}^*, \text{let } x = v \text{ in } t \rightarrow \mathcal{V}_0[x \leftarrow v] \mathcal{V}^*, t}$$

A let-binding $\text{let } x = v \text{ in } t$ reduces to t and adds a binding from x to v in the current environment.

$$\text{application} \frac{\begin{array}{l} \mathcal{V}_0(x_1) = a \quad \mathcal{M}(a) = \mathcal{S}.x \\ \mathcal{S} = _ \mathcal{C} \quad \mathcal{C}(x) = \lambda x_3 x_4 . t \\ \mathcal{V}_0(x_2) = v \end{array}}{\mathcal{V}_0 \mathcal{V}^*, \mathcal{M}, x_1 x_2 \rightarrow [x_3 \leftarrow v; x_4 \leftarrow a] \mathcal{V}_0 \mathcal{V}^*, \text{pop } t}$$

An application $x_1 x_2$ looks for an address a in the current environment. The image of a in the store must be a closure $\mathcal{S}.x$ where the set of closures $\mathcal{S} = _ \mathcal{C}$. Let the value of x in the closure map \mathcal{C} be $\lambda x_3 x_4 . t$. The term reduces to $\text{pop } t$ and we create a new environment containing two bindings: x_3 is mapped to the value of x_2 in the current environment and x_4 is mapped to a .

365

366

367

$$\text{pop environment} \frac{}{\mathcal{V}_0 \mathcal{V}^*, \text{pop } v \rightarrow \mathcal{V}^*, v}$$

368

369

This simply throws away the current environment. This is used to model the call stack. A pop is inserted around each function application and removed once the function is fully evaluated to a value.

370

371

372

373

374

$$\text{mutate} \frac{\mathcal{V}_0(x_1) = v_1 \quad \mathcal{V}_0(x_2) = v_2}{\mathcal{V}_0 \mathcal{V}^*, x_1 <- x_2 \rightarrow \mathcal{V}_0[x_1 \leftarrow v_2] \mathcal{V}^*, 0}$$

375

376

377

A mutation $x_1 <- x_2$ simply updates the value of x_1 to be the value of x_2 in the current environment. It reduces to 0 which corresponds to the value `()` : **unit** in OCaml.

378

379

380

381

382

383

$$\text{conditional 1} \frac{\mathcal{V}_0(x) = n \quad n \neq 0}{\mathcal{V}_0 \mathcal{V}^*, \text{if } x \text{ then } t_1 \text{ else } t_2} \rightarrow t_1$$

384

385

386

387

388

389

390

391

A conditional `if x then t_1 else t_2` reduces to t_1 when the value of x in the current environment is not 0 (i.e. when it is **true** : **bool** in OCaml).

392

393

394

395

396

$$\text{conditional 2} \frac{\mathcal{V}_0(x) = n \quad n = 0}{\mathcal{V}_0 \mathcal{V}^*, \text{if } x \text{ then } t_1 \text{ else } t_2} \rightarrow t_2$$

Similarly, it reduces to t_2 when the value of x is 0 (i.e. when it is **false** : **bool** in OCaml).

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

When the value of x in the current environment is an integer n , then `switch x \mathcal{B}_1 \mathcal{B}_2` reduces to the image of n in \mathcal{B}_1 .

$$\text{switch 1} \frac{\mathcal{V}_0(x) = n \quad \mathcal{B}_1(n) = t}{\mathcal{V}_0 \mathcal{V}^*, \text{switch } x \mathcal{B}_1 \mathcal{B}_2 \rightarrow t}$$

Similarly, when x is an address a and $\mathcal{M}(a)$ is a block with tag n , then `switch x \mathcal{B}_1 \mathcal{B}_2` reduces to the image of n in \mathcal{B}_2 .

$$\text{switch 2} \frac{\mathcal{V}_0(x) = a \quad \mathcal{M}(a) = n, v^* \quad \mathcal{B}_2(n) = t}{\mathcal{V}_0 \mathcal{V}^*, \mathcal{M}, \text{switch } x \mathcal{B}_1 \mathcal{B}_2 \rightarrow t}$$

If a raised exception does not match the exception expected by a static catch, then the catch is discarded and we re-raise the exception.

$$\text{raise 2} \frac{\begin{array}{l} x_{exn} = x'_{exn} \quad x^n = x_0, \dots, x_{n-1} \\ x'^n = x'_0, \dots, x'_{n-1} \quad \forall i \in [0; n], \mathcal{V}'_0(x_i) = v_i \end{array}}{\begin{array}{l} \mathcal{V}'_0 \mathcal{V}^*, \\ \text{scatch (sraise } x_{exn} x^n) \text{ with } x'_{exn} x'^n t \\ \rightarrow \mathcal{V}'_0[x'_0 \leftarrow v_0, \dots, x'_{n-1} \leftarrow v_{n-1}] \mathcal{V}^*, \\ t \end{array}}$$

On the contrary, when the two exceptions match, the identifiers specified in the catch are binded to the values of the identifiers transported by the exception. The values are taken from the current environment. We have the guarantee that if the transported identifiers were in the environment when we raised, then they are also in the environment when we catch. This is because static exceptions never cross a function application. Finally, the reduced term is the one attached to the catch.

$$\text{raise 3} \frac{\text{while}_{t'} (\text{sraise } x_{exn} x^*) \text{ do } t}{\rightarrow \text{sraise } x_{exn} x^*}$$

A while simply re-raises static exceptions.

$$\text{raise 4} \frac{\text{let } x = (\text{sraise } x_{exn} x^*) \text{ in } t}{\rightarrow \text{sraise } x_{exn} x^*}$$

Similarly, a let-binding also re-raises static exceptions.

$$\text{while loop 1} \frac{n = 0}{\text{while}_{t'} n \text{ do } t \rightarrow 0}$$

A loop while_{t'} n do t reduces to 0 when n = 0.

$$\text{while loop 2} \frac{n \neq 0}{\begin{array}{l} \text{while}_{t'} n \text{ do } t \\ \rightarrow \text{let } _ = t \text{ in while}_{t'} t' \text{ do } t \end{array}}$$

On the contrary, when n ≠ 0, it reduces to let ₋ = t in while_{t'} t' do t. The term t' corresponds to the original term that was used as a condition.

$$\text{field read} \frac{\begin{array}{l} \mathcal{V}'_0(x) = a \quad \mathcal{M}(a) = \text{tag}, v^* \\ v^* = v_0, \dots, v_{m-1} \quad m > n \geq 0 \end{array}}{\mathcal{V}'_0 \mathcal{V}^*, \mathcal{M}, \text{get_field } n \ x \rightarrow v_n}$$

The term get_field n x reduces to the nth value of the block stored at address a with a being the value of x in the current environment.

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

$$\begin{array}{c}
 \mathcal{V}_0(x_1) = a \quad \mathcal{M}(a) = \text{tag}, v^* \quad v^* = v_0, \dots, v_{m-1} \\
 m > n \geq 0 \quad \mathcal{V}_0(x_2) = v \\
 \text{field write} \text{---} \\
 \mathcal{V}_0 \mathcal{V}^*, \\
 \mathcal{M}, \\
 \text{set_field } n \ x_1 \ x_2 \\
 \rightarrow \mathcal{M}[a \leftarrow \text{tag}, v_0, \dots, v_{n-1}, v, v_{n+1}, \dots, v_{m-1}], \\
 0
 \end{array}$$

If x_1 is the address a in the current environment, then the term `set_field n x_1 x_2` sets the n^{th} value of the block stored at address a to the value of x_2 in the current environment.

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

$$\begin{array}{c}
 \mathcal{V}_0(x_0) = v_0, \dots, \mathcal{V}_0(x_{n-1}) = v_{n-1} \quad a \notin \text{dom}(\mathcal{M}) \\
 \text{block creation} \text{---} \\
 \mathcal{V}_0 \mathcal{V}^*, \mathcal{M}, \text{make_block } n_{\text{tag}} \ x_0, \dots, x_{n-1} \\
 \rightarrow \mathcal{M}[a \leftarrow n_{\text{tag}}, v_0, \dots, v_{n-1}], 0
 \end{array}$$

This creates a new block in the store at a fresh address. The tag is given as a value while the others values of the block are read from the current environment.

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

$$\begin{array}{c}
 \mathcal{V}_0(x_2) = a_1 \\
 \mathcal{M}(a_1) = \mathcal{S} \quad \mathcal{M}(a_2) = \mathcal{S}.x_1 \\
 \text{closure projection} \text{---} \\
 \mathcal{V}_0 \mathcal{V}^*, \mathcal{M}, \text{project_closure } x_1 \ x_2 \\
 \rightarrow a_2
 \end{array}$$

This reduces to the address a_2 of the closure which has the identifier x_1 in the set of closures at address a_1 , with a_1 being the value of x_2 in the current environment. In short, it allows one to get a particular function from a previously allocated set of closures.

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

$$\begin{array}{c}
 \mathcal{V}_0(x_2) = a \quad \mathcal{M}(a) = \mathcal{S}.x \\
 \mathcal{S} = \mathcal{V}, _ \quad \mathcal{V}(x_1) = v \\
 \text{variable projection} \text{---} \\
 \mathcal{V}_0 \mathcal{V}^*, \mathcal{M}, \text{project_var } x_1 \ x_2 \rightarrow v
 \end{array}$$

This reduces to the value of x_1 in the environment of the set of closures that contains the closure at address a_a , with a_a being the value of x_2 in the current environment. In short, when inside a closure, it allows to read a variable from the set of closures the closure belongs to.

510

511

512

513

514

515

516

517

518

519

520

$$\begin{array}{c}
 \mathcal{V}_0(x_2) = a_1 \\
 \mathcal{M}(a_1) = \mathcal{S}.x \quad \mathcal{M}(a_2) = \mathcal{S}.x_1 \\
 \text{closure movement} \text{---} \\
 \mathcal{V}_0 \mathcal{V}^*, \mathcal{M}, \text{move_within_closure } x_1 \ x_2 \\
 \rightarrow a_2
 \end{array}$$

This reduces to the address a_2 of the closure x_1 in the set of closures to which the closure at address a_1 belongs to. In short, when inside a closure, it allows to get another closure in the current set of closures.

$$\begin{array}{c}
521 \\
522 \\
523 \\
524 \\
525 \\
526 \\
527 \\
528
\end{array}
\begin{array}{c}
\mathcal{S} = _, \mathcal{C} \quad x_0 \in \text{dom}(\mathcal{C}), \dots, x_{n-1} \in \text{dom}(\mathcal{C}) \\
\forall a \in [a_0, \dots, a_n], a \notin \text{dom}(\mathcal{M}) \\
\hline
\text{set alloc} \quad \mathcal{M}, \text{make_set_of_closures } \mathcal{S} \\
\rightarrow \mathcal{M} [a_n \leftarrow \mathcal{S}; a_0 \leftarrow \mathcal{S}.x_0; \dots; a_{n-1} \leftarrow \mathcal{S}.x_{n-1}], \\
a_n
\end{array}$$

529 This maps the set of closures \mathcal{S} and all its closures to fresh addresses in the store. The term reduces to the address of
530 the set.

532 2.4 Example

533 Consider the following OCaml code:

```
535 let f = print_int
536 let rec iter f l =
537   match l with
538   | [] -> ()
539   | hd :: tl ->
540     f hd;
541     iter f tl
```

544 In Flambda, it would be represented as follows:

```
545
546 let f = (* ... *)
547 let iter =
548   let set =
549     make_closures
550     | cl_iter { f } env ->
551       let otherset =
552         make_set_of_closures
553         (* this is the closure map *)
554         | cl_iter_f { l } envtwo ->
555           switch l
556           | int 0 -> 0
557           | tag 0 ->
558             let x = get_field 0 l in
559             let f = project_var f envtwo in
560             let dummy = f x in
561             let tl = get_field 1 l in
562             envtwo tl
563           (* this is the environment *)
564           | f -> f
565         end
566   in
```

```

573     project_closure cl_iter_f otherset
574   in
575     project_closure cl_iter set
576   in
577     (* ... *)
578
579

```

3 Compiling Flambda to Wasm

Wasm code runs inside an *embedder*. In the case of the web it is the browser and the host language is JavaScript, which has a GC. As OCaml also has a GC, it means having two different runtimes and two different garbage collectors. For this reason, it is hard to avoid memory leaks when compiling OCaml to Wasm1. Indeed, cycles between the two GCs can not be collected without adding a third GC. it is hard to avoid memory leaks when compiling OCaml to Wasm1, because it means having two different runtimes and two different garbage collectors.

To properly interact with the embedder, we need to leverage the reference and the GC extensions. This extension adds new types to the language, e.g. `externref` which is an opaque type representing a value from the embedder. References cannot be stored in the linear memory of Wasm thus they cannot appear inside OCaml values when using the previously described compilation scheme.

In order to use references, we require a completely different compilation strategy. We do not use the linear memory but rely entirely on WasmGC. Our strategy is close to the native OCaml one, which we describe now.

We use the Flambda IR of the OCaml compiler as input for the Wasm generation. This is a step of the compilation chain where most of the high-level OCaml-specific optimisations are already applied. Also in this IR, the closure conversion pass is already performed. Most of the constructions of this IR maps quite directly to Wasm ones.

Wasocaml is a compiler for the full Flambda IR. We describe the full compilation scheme but for the sake of simplicity we only formalize a subset of Flambda.

3.1 OCaml Value Representation in Wasm

As in native OCaml, we use a uniform representation. We cannot use integers, as they cannot be used as pointers. We need to use a reference type. The most general one is `anyref`. We can be more precise and use `eqref`. This is the type of all values that can be tested for physical equality. It is a super type of OCaml values. This also allows us to test for equality without requiring an additional cast.

Small scalars. All OCaml values that are represented as integers in the native OCaml are represented as `i31ref`. This includes the `int` and `char` types, but also all constant constructors of sum types.

Arrays. The OCaml array type is directly represented as a Wasm array.

Blocks. For other kinds of blocks, there are two possible choices: we can either use a struct with a field for the tag and one field per OCaml value field, or use arrays of `eqref` with the tag stored at position 0.

3.1.1 Blocks as Structs. A block of size one is represented using the type `$block1` while for size two it is `$block2`:

```

620 (type $block1 (struct
621   (field $tag i8)
622   (field $field0 eqref)))

```

```

625
626 (type $block2 (sub $block1) (struct
627   (field $tag i8)
628   (field $field0 eqref)
629   (field $field1 eqref)))
630
631

```

The type `$block2` is declared as a sub-type of `$block1` because in the OCaml IRs, the primitive for accessing a block field is untyped: when accessing the n -th field of the block the only available information is that the block has size at least $n + 1$. It could be possible to propagate size metadata in the IRs but that wouldn't be sufficient because of untyped primitives such as `Obj.field` that we need to support in order to be compatible with existing code.

```

637
638 (func $snd (param $x eqref) (result eqref)
639   (struct.get $block2 1
640     (ref.cast $block2 (local.get $x))))
641
642

```

3.1.2 *Blocks as Arrays.* We represent blocks with an array of eqref:

```

644 (type $block (array eqref))
645
646

```

The tag is stored in the cell at index 0 of the array. Reading its value is implemented by getting the cell and casting it to an integer:

```

649
650 (func $tag (param $x eqref) (result i32)
651   (i31.get_u (ref.cast i31 ref
652     (array.get $block
653       (i32.const 0)
654       (ref.cast $block (local.get $x)))))
655
656

```

Thus accessing the field n of the OCaml block amounts to accessing the field $n + 1$ of the array:

```

658 (func $snd (param $x eqref) (result eqref)
659   (array.get $block
660     (i32.const 2)
661     (ref.cast $block (local.get $x))))
662
663
664

```

3.1.3 *Block Representation Tradeoffs.* The array representation is simpler but requires (implicit) bound checks at each field access and a cast to read the tag.

On the other hand, the struct representation requires a more complex cast for the Wasm runtime (a sub-typing test). A compiler propagating more types could use finer Wasm type information, providing a precise type for each struct field. This would allow fewer casts.

OCaml blocks can be arbitrarily large, and very large ones do occur in practice. In particular, modules are compiled as blocks and tend to be on the larger side. We have seen in the wild some examples containing thousands of fields. Wasm only allows a sub-typing chain of length 64, which is far from sufficient for the struct encoding.

For this reason we use the *blocks as arrays* representation.

677 3.1.4 *Boxed Numbers*. Raw Wasm scalars such as `i64` are not sub-types of `eqref` thus they cannot be used directly
 678 to represent OCaml boxed numbers. We need to box them inside a struct in order to make them compatible with our
 679 representation:
 680

```
681 (type $boxed_float (struct (field $v f64)))
682 (type $boxed_int64 (struct (field $v i64)))
```

684 3.1.5 *Closures*. Wasm `funcref` are functions, not closures, hence we need to produce values containing both the
 685 function and its environment. The only Wasm type construction that can contain both `funcref` and other values are
 686 structs. Thus, a closure is a struct containing a `funcref` and the captured variables. As an example, here is the type of
 687 a closure with two captured variables:
 688

```
689
690 (type $closure1 (struct
691   (field funcref)
692   (field $v1 eqref)
693   (field $v2 eqref)))
```

696 3.2 Control flow

697 Control flow has a direct equivalent with Wasm `block`, `loop`, `br_table`, and `if` instructions. Low level OCaml prim-
 698 itives to handle exceptions are quite similar to Wasm ones. In OCaml, it is possible to generate new exceptions at
 699 runtime by using e.g. the `let exception` syntax or functors and first-class modules. This is not possible in the Wasm
 700 exception proposal. Thus, we use the same Wasm exception everywhere and manage the rest on the side by ourselves,
 701 using an identifier to discriminate between different exceptions.
 702
 703

705 3.3 Curryfication

706 In OCaml, functions take only one argument. However, in practice, functions look like they have more than one. Without
 707 any special management this would mean that most of the code would be functions producing closures that would
 708 be immediately applied. To handle that, internally, OCaml does have functions taking multiple arguments, with some
 709 special handling for times when they are partially applied. This information is explicit at the Lambda level. In the native
 710 OCaml compiler, the transformation handling that occurs in a later step called `cmmgen`. Hence, we have to duplicate this
 711 in Wasocaml. Compiling this requires some kind of structural sub-typing on closures such that, closures for functions
 712 of arity one are supertypes of all the other closures. Thankfully there are easy encodings for that in Wasm.
 713
 714

716 3.4 Stack Representation

718 Most of the remaining work revolves around translating from let bound expressions to a stack based language without
 719 producing overly naive code. Also, we do not need to care too much about low level Wasm specific optimisations as
 720 we rely on Binaryen [16] (a quite efficient Wasm to Wasm optimizer) for those.
 721

723 3.5 Unboxing

724 The main optimisation available in native OCaml that we are missing is number unboxing. As OCaml values have a
 725 uniform representation, only small scalars can fit in a true OCaml value. This means that the types `nativeint`, `int32`,
 726 `int64` and `float` have to be boxed. In numerical code, lots of intermediate values of type float are created, and in that

case, the allocation time to box numbers completely dominates the actual computation time. To limit that problem, there is an optimisation called unboxing performed during the `cmmgen` pass that tends to eliminate most of the useless allocations. As this pass is performed after `Flambda`, and was not required to produce a complete working compiler, this was left for future work. Note that the end plan is to use the next version of `Flambda`, which does a much better job at unboxing. For the time being, we can expect `Binaryen` to perform local unboxing in some cases.

4 Formalized compilation scheme

When compiling a program, we start with the following prelude:

```
(module $wasocaml_generated_modul

  (type $mlfun (sub final
    (func (param eqref) (param eqref) (result eqref))))

  (type $data (sub final
    (array (mut (ref null eq))))))

  (type $block (sub final (struct
    (field i32)
    (field (ref null $data)))))

  (type $vars (sub final (array (mut eqref))))

  (rec
    (type $set_of_closures (sub final (struct
      (field (mut (ref null $closures)))
      (field (ref null $vars)))))

    (type $closures (sub final
      (array (mut (ref $closure)))))

    (type $closure (sub final (struct
      (field (ref $mlfun))
      (field (mut (ref null $set_of_closures)))))))

  (global $tmp_set_of_closures
    (mut (ref null $set_of_closures))
    ref.null $set_of_closures)
  (global $switch_v
    (mut i32) i32.const 0)
```



```

781 (elem declare (ref $m1fun)
782   ;; to allow forward references of functions
783   ;; ...
784 )
785
786
787 (func $start
788   ;; ...
789   drop
790 )
791
792
793 (start $start))
794
795
796
797
798
799
800
801
802
803
804

```

805 To compile a term we define a function \mathcal{T} from Flambda terms to Wasm expressions. We define it case by case.
 806 Here we do not only use the S-expression notation for Wasm but also use the stack notation which makes these rules
 807 clearer.

$$810 \mathcal{T}(n) = i32.const\ n$$

811 `i31.new`

812 An integer literal is turned into an `i31.ref` in order to match with the uniform representation.

$$817 \mathcal{T}(x) = local.get\ $x$$

818 To access the value of an identifier, we simply access the value stored in the local variable of the current function.

$$822 \mathcal{T}(\text{let } x = t_1 \text{ in } t_2) = \mathcal{T}(t_1)$$

823 `local.set $x`

824 $\mathcal{T}(t_2)$

825 A let-binding is compiled by evaluating t_1 , storing its result in the local x and evaluating t_2 .

```

833
834
835  $\mathcal{T}(x_1\ x_2) = \text{local.get } \$x_1$ 
836  $\text{ref.cast (ref } \$\text{closure)}$ 
837
838  $\text{local.get } \$x_2$ 
839  $\text{ref.as\_non\_null}$ 
840
841  $\text{local.get } \$x_1$ 
842  $\text{ref.cast (ref } \$\text{closure)}$ 
843
844  $\text{struct.get } \$\text{closure } 0$ 
845
846  $\text{call\_ref } \$\text{mlfun}$ 
847

```

To apply a function x_1 to a value x_2 we first put the two parameters of the function on the stack. The first one is the closure itself, which is needed to access the environment and other functions in the set of closures. The second one can be any value of type `eqref`. Then we need to put the function of the stack. This is done by accessing the field 0 of the closure, which is a `funcref`. Then we use `call_ref` to call the function with its two arguments. The type of the function reference is always `$mlfun`.

```

855
856
857  $\mathcal{T}(x_1 <- x_2) = \text{local.get } \$x_2$ 
858  $\text{local.set } \$x_1$ 
859
860  $\text{i32.const } 0$ 
861
862  $\text{i31.new}$ 
863

```

A mutation simply puts the content of x_2 into x_1 and leaves 0 on the stack (this is `()` : `unit` in OCaml).

```

865
866
867  $\mathcal{T}(\text{if } x_1 \text{ then } t_1 \text{ else } t_2) = \text{local.get } \$x_1$ 
868  $\text{ref.cast } \text{i31.ref}$ 
869
870  $\text{i31.get\_s}$ 
871
872  $(\text{if (result (ref null eq))$ 
873  $\text{(then } \mathcal{T}(t_1))$ 
874  $\text{(else } \mathcal{T}(t_2))$ 
875
876  $\text{ref.as\_non\_null}$ 
877

```

A conditional converts the value of x_1 to an `i32` and then uses a Wasm `if`. The two branches are simply the result of compiling t_1 and t_2 .

```

881
882
883
884

```

```

885
886
887  $\mathcal{T}(\text{switch } x \mathcal{B}_1 \mathcal{B}_2) =$ 
888
889 (block $switch_a (result (ref null eq))
890
891   (block $switch_b (result (ref null eq))
892
893     local.get $x
894     br_on_cast $switch_b (ref null eq) i31.ref
895     ;; handle block case
896     ref.cast (ref $block)
897     struct.get $block 0
898     global.set $switch_v
899      $\mathcal{F}(\mathcal{B}_2)$ 
900     ;; skip the scalar case
901     br $switch_a)
902
903   ;; handle scalar case
904   ref.cast (ref i31)
905   i31.get_s
906   global.set $switch_v
907    $\mathcal{F}(\mathcal{B}_1)$ )
908
909
910
911
912
913
914
915

```

To compile a switch, we start by testing if x is a block or a scalar value. If it's a block, we store its tag, otherwise we store the value directly. Then, we use an auxiliary function \mathcal{F} in order to compile the various cases of each branches. Let $i_0, \dots, i_{n-1} \in \text{dom}(\mathcal{B})$ be such that $\nexists j, j \notin \text{dom}(\mathcal{B}) \wedge \forall k \in [0; n[, j \neq i_k$. Then we define the auxiliary function as follows:

```

921
922  $\mathcal{F}(\mathcal{B}) =$ 
923
924 (i32.eq (i32.const  $i_{n-1}$ ) (global.get $switch_v))
925
926 (if (result (ref null eq)) (then  $\mathcal{F}(\mathcal{B}(i_{n-1}))$ )
927
928 (else
929
930   ;; ...
931   (i32.eq (i32.const  $i_0$ ) (global.get $switch_v))
932   (if (result (ref null eq)) (then  $\mathcal{F}(\mathcal{B}(i_0))$ )
933   (else unreachable))))
934
935
936

```

We have a sequence of nested `if` to find which of the binding in \mathcal{B} we should evaluate. If none of them correspond to the value stored in `$switch_v` (which has been set by the code generated when compiling the switch) then we have to put an `unreachable` in order to generate Wasm code that typechecks. In practice, the `unreachable` corresponds to a match failure, which should not happen in a valid source program, thus it is indeed unreachable.

```

937  $\mathcal{T}(\text{srise } x_{\text{exn}} x^n) = \text{local.get } \$x_{n-1}$ 
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988

```

A static raise puts all the value transported by the exception on the stack and branches to the block corresponding to the raise.

```

937  $\mathcal{T}(\text{scatch } t_1 \text{ with } x_{\text{exn}} x^n t_2) =$ 
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988

```

A catch has two nested blocks. The first one is used in case no exception is raised and directly returns the value of t_1 . The second one corresponds to the case where an exception is raised. It needs to have many results corresponding to all the values transported by the exception. It will put all these values in the right locals.

```

989
990
991  $\mathcal{T}(\text{while}_{t_1} t_1 \text{ do } t_2) = (\text{block } \$\text{exit\_while\_loop}$ 
992      $(\text{loop } \$\text{while\_loop}$ 
993          $\mathcal{T}(t_1)$ 
994          $i31.\text{get\_s}$ 
995          $i32.\text{eqz}$ 
996          $\text{br\_if } \$\text{exit\_while\_loop}$ 
997          $\mathcal{T}(t_2)$ 
998          $\text{drop}$ 
999          $\text{br } \$\text{while\_loop} ) )$ 
1000
1001
1002
1003
1004
1005      $i32.\text{const } 0$ 
1006      $i31.\text{new}$ 

```

A while loop maps quite directly to a Wasm **loop**. We also need an other block in order to exit the loop when the condition is false.

```

1011
1012
1013  $\mathcal{T}(\text{get\_field } n \ x) = \text{local.get } \$x$ 
1014      $\text{ref.cast } (\text{ref } \$\text{block})$ 
1015      $\text{struct.get } \$\text{block } 1$ 
1016      $\text{ref.as\_non\_null}$ 
1017      $i32.\text{const } n$ 
1018      $\text{array.get } \$\text{data}$ 

```

To read a field of a block, we need to cast it to a block, access its second field which contains the values of the block (the first field being the tag). Then we simply read the element at the index n .

```

1041
1042
1043  $\mathcal{T}(\text{set\_field } n \ x_1 \ x_2) = \text{local.get } \$x_1$ 
1044  $\text{ref.cast (ref } \$\text{block)}$ 
1045  $\text{struct.get } \$\text{block } 1$ 
1046  $\text{ref.cast (ref } \$\text{data)}$ 
1047  $\text{i32.const } x$ 
1048  $\text{local.get } \$x_2$ 
1049  $\text{array.set } \$\text{data}$ 
1050  $\text{i32.const } 0$ 
1051  $\text{i31.new}$ 

```

Similarly, to set a field of a block, we cast it, access its data field and then set the element at index n to the value of x_2 . The array is mutable and there is no need to update the structure after the element has been set.

```

1061
1062  $\mathcal{T}(\text{make\_block } n \ x_0, \dots, x_{m-1}) = \text{i32.const } n$ 
1063  $\text{local.get } \$x_0$ 
1064  $\dots$ 
1065  $\text{local.get } \$x_{m-1}$ 
1066  $\text{array.new\_fixed } \$\text{data } m$ 
1067  $\text{struct.new } \$\text{block}$ 

```

To create a new block, we put its tag on the stack. Then we put all its data and make a new array from it. With the tag and the array we then make a new structure of type $\$block$.

```

1072
1073  $\mathcal{T}(\text{project\_closure } x_1 \ x_2) =$ 
1074  $\text{local.get } \$x_2$ 
1075  $\text{ref.cast (ref } \$\text{set\_of\_closures)}$ 
1076  $\text{struct.get } \$\text{set\_of\_closures } 0$ 
1077  $\text{global.get } \$\text{closure\_offset\_within\_set\_}x_1$ 
1078  $\text{array.get } \$\text{closures}$ 

```

To read a closure from a set of closures, we cast x_2 to a set. Then we access its first field, which is an array of closures. Then we look for the offset of x_1 and access it.

```

1089
1090
1091
1092

```

```

1093
1094
1095  $\mathcal{T}(\text{project\_var } x_1 \ x_2) =$ 
1096 local.get $x_2
1097
1098 ref.cast (ref $closure)
1099
1100 struct.get $closure 1
1101
1102 ref.as_non_null
1103
1104 struct.get $set_of_closures 1
1105
1106 global.get $var_offset_within_set_x1
1107
1108 array.get $vars

```

This is similar to the closure case but the difference here is that x_2 is a closure and not a set. We need to access its parent first. Then we access the second field, which is the array of the captured variables. Then we need to look for the offset of the variable again before reading the array.

```

1113
1114  $\mathcal{T}(\text{move\_within\_closure } x_1 \ x_2) =$ 
1115 local.get $x_2
1116
1117 ref.cast (ref $closure)
1118
1119 struct.get $closure 1
1120
1121 ref.as_non_null
1122
1123 struct.get $set_of_closures 0
1124
1125 global.get $closure_offset_within_set_x1
1126
1127 array.get $closures

```

This is the same as the previous case: x_2 is a closure and not a set. The only difference being that we look in the first field instead of the second one.

1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144

For allocating a set of closures $\mathcal{S} = \mathcal{V}, \mathcal{C}$ we define x_0, \dots, x_{n-1} to be the values that are captured in \mathcal{V} and f_0, \dots, f_{m-1} to be the names of the functions that are part of the set.

```

1145      $\mathcal{T}(\text{make\_set\_of\_closures } \mathcal{S}) =$ 
1146     ;; empty closures
1147
1148     ref.null $closures
1149
1150     ;; generating variables
1151     (ref.as_non_null (local.get $x0))
1152
1153     ;; ...
1154     (ref.as_non_null (local.get $xn-1))
1155
1156     array.new_fixed $vars n
1157
1158     ;; set of closures with empty closures and filled vars
1159     struct.new $set_of_closures
1160
1161     global.set $tmp_set_of_closures
1162     global.get $tmp_set_of_closures
1163
1164     ;; generating closures
1165     ref.func $f0
1166     global.get $tmp_set_of_closures
1167     struct.new $closure
1168
1169     ;; ...
1170     ref.func $fm-1
1171     global.get $tmp_set_of_closures
1172     struct.new $closure
1173
1174     array.new_fixed $closures m
1175
1176     ;; patching the set with closures
1177     struct.set $set_of_closures 0
1178
1179     ;; returning the set
1180     global.get $tmp_set_of_closures
1181
1182
1183
1184
1185
1186
1187
1188

```

We start by generating an empty array of closures. We build the array made of the captured variables. Then we generate an empty set of closures from our empty array of closures and our array of variables.. Each closure needs a reference to its set, hence we have to store the set in a global before constructing them. Then, we build the array made of the various functions in the set and patch the set.

1194
1195
1196

1197 The code of the compiled Flambda program is inserted as the code of the `$start` function. But we also need to
 1198 analyse the term in order to do a few more things. Each function defined in a set of closures is also added as a new
 1199 Wasm function and marked as a declarative element (in order to be referenced by the `ref.func` instruction). More-
 1200 over, each function and each variable in a set of closures has its own offset stored as a global variable (for instance
 1201 `global $closure_offset_within_set_myfunc`). Finally, each function is analysed in order to find which locals it
 1202 requires to be declared.
 1203

1204 For now, we don't have a proof of the correctness of our compilation scheme. The formal semantics of the various
 1205 Wasm extensions we are needing is still an ongoing work by the various working groups. We would like to use these
 1206 semantics in the future.
 1207
 1208

1209 4.1 Example of generated code

1210 The code of `cl_iter_f` from the example given in 2.4 would be compiled to the following Wasm code:

```

1211
1212 (func $cl_iter_f (param $envtwo eqref) (param $l eqref)
1213 (result eqref) (local $x eqref) (local $f eqref)
1214 (local $dummy eqref) (local $t1 eqref)
1215 (block $switch_a (result (ref null eq))
1216 (block $switch_b (result (ref null eq))
1217 local.get $l
1218 br_on_cast $switch_b (ref null eq) i31.ref
1219 ref.cast (ref $block)
1220 struct.get $block 0
1221 global.set $switch_v
1222 global.get $switch_v
1223 i32.const 0
1224 i32.eq
1225 (if (result (ref null eq))
1226 (then
1227 local.get $l
1228 ref.cast (ref $block)
1229 struct.get $block 1
1230 ;; ...
1231 local.get $envtwo
1232 ref.cast (ref $closure)
1233 struct.get $closure 0
1234 call_ref $mlfun)
1235 (else unreachable))
1236 br $switch_a)
1237 ref.cast (ref i31)
1238 i31.get_s
1239 global.set $switch_v
1240

```

```
1249     global.get $switch_v
1250     i32.const 0
1251     i32.eq
1252     (if (result (ref null eq))
1253       (then
1254         i32.const 0
1255         i31.new)
1256       (else
1257         unreachable))))
```

5 Performances

We are able to run the classical functional microbenchmarks using the experimental branch of V8 and the results are quite convincing. While we do not expect to be competitive with the native code compiler, the performance degradation is almost constant (around twice times slower).

We are also able to compile an OCaml implementation of the Knuth-Bendix algorithm [5]. On V8 the Wasm exception handling proposal is really slow compared to native OCaml: a raise is around a hundred times slower. In SpiderMonkey (the Firefox engine) they are much faster but other proposals are missing and we can't run code we generate. We have to discuss with the V8 team if this is expected and if this can be improved. We implemented an alternative compilation strategy for exceptions where we do not use Wasm exceptions. Instead we make every function return a pair containing its return value along with a boolean indicating if an exception was raised or not. When using this alternative strategy, we are around two times slower than native OCaml.

The V8 runtime allows to consider casts as no-ops (only for test purposes). The speedup we measured is around 10%. That gives us an upper bound on the actual runtime cost of casts.

We are currently working on producing benchmarks on real sized programs, it is not easy as Wasocaml is still a prototype and is not yet integrated into existing build systems.

Compared to a JavaScript VM, a Wasm compiler is a much simpler beast that can compile ahead of time. For this reason, various Wasm engines are expected to behave quite similarly. They do not show any of the wild unpredictability that browsers tend to demonstrate with JavaScript. Indeed, compiling OCaml to JS using `js_of_ocaml` leads to results that are usually also twice as slow as native code in the best cases, but can sometimes be much slower in an unpredictable fashion.

Currently there is no other Wasm runtime supporting all the extensions we require. SpiderMonkey does not have tail calls. The reference interpreter implementation of the various extensions are split in separate repositories and merging them requires some work.

6 Perspectives

As the first version of the implementation was intended as a demonstration, it remains a bit rough around the edges. In particular only a fraction of the externals from the standard library externals are provided as handwritten Wasm. The only unsupported part of the language is the object-oriented fragment, by lack of time rather than due to any specific complexity: in Flambda this fragment is reduced to a single primitive. The source code of Wasocaml is publicly available [1].

1301 6.1 FFI

1302 A lot of OCaml code in the wild interacts with C or JavaScript code through bindings written using dedicated FFI
 1303 mechanisms. We plan to allow re-use of existing bindings when compiling to Wasm. When targeting the browser
 1304 through JavaScript, it is not possible to re-use C bindings: people either have to rewrite their code in pure OCaml, or
 1305 to use JavaScript bindings. Compiling C bindings directly to Wasm allows to re-use them as-is.
 1306
 1307

1308 *C bindings.* Some extensions recently added to Clang [6] allow to compile C code to Wasm in a way that makes
 1309 reusing existing OCaml bindings to C code possible with almost no change. We would only have to provide a modified
 1310 version of the OCaml native FFI headers files, replacing the usual macros with hand written Wasm functions. The only
 1311 limitation that we foresee is that the `Field` macro won't be an l-value anymore; a new `Set_field` macro will be needed
 1312 instead.
 1313
 1314

1315 *JavaScript bindings.* To re-use existing bindings of OCaml to JavaScript code, we need one more extension: the
 1316 reference-typed strings proposal [24]. Almost all external calls in the JavaScript FFI go through the `Js.Unsafe.meth_call`
 1317 function which has the type `'a -> string -> any array -> 'b`. It can be exposed to the Wasm module by the em-
 1318 bedder as a function of type:
 1319
 1320

```
1321 (func $meth_call
1322   (param $obj externref)
1323   (param $method stringref)
1324   (param $args $anyarray)
1325   (result externref))
```

1326
 1327
 1328 This calls a method of name `$method` on the object `$obj` with the arguments `$args`. The JavaScript side manages
 1329 all the dynamic typing.
 1330

1331 6.2 Effect Handlers Support

1332
 1333 Our compiler is based on OCaml 4.14. The OCaml 5 compiler introduces effect handlers [10], a mechanism that can be
 1334 seen as a generalisation of exceptions, allowing them to be resumed. They were out of scope for our compiler, yet we
 1335 describe three strategies to handle them in the future.
 1336
 1337

1338 *CPS Compilation.* It is possible to represent effect handlers as continuations, using whole-program CPS transforma-
 1339 tion [4]. However it is not ideal for two reasons. It requires the program to contain OCaml code only. Moreover, it
 1340 prevents the use of the stack of the target language and instead stores that stack in the heap, with a non-negligible cost.
 1341 It can be mitigated by only applying the transformation to code that cannot be proved to not use effect handlers but
 1342 then it is no longer compatible with separate compilation and the optimisation needs to be done as a whole program
 1343 transformation. This is the choice made by `Js_of_ocaml`.
 1344
 1345

1346 *Wasm Stack Switching.* There is an ongoing proposal, called stack switching [23], that exactly matches the needs for
 1347 OCaml effect handlers. With it the translation would be quite straightforward.
 1348
 1349

1350 *JavaScript Promise Integration.* Another strategy is available for runtimes providing both Wasm and JavaScript.
 1351 There is another ongoing proposal called JavaScript promise integration [22]. With this proposal, effects handlers
 1352

1353 can be implemented with a JavaScript device. This proposal is likely to land before the proper stack switching one. It
1354 might be a temporary solution.
1355

1356 7 Related Work

1357 7.1 Garbage Collected Languages to Wasm

1358 7.1.1 Targeting Wasm1.

1359
1360
1361 *Go*. Go compiles to Wasm1 [9]. In order to be able to re-use the Go GC, it must be able to inspect the stack. This is
1362 not possible within Wasm. Thus it has to maintain the stack by itself using the memory, which is much slower than
1363 using the Wasm stack.
1364
1365

1366 *Haskell*. Haskell also compiles to Wasm [11]. It has constraints similar to Go and uses similar solutions. In particular,
1367 it also uses the memory to manage the stack.
1368

1369 7.1.2 Targeting Wasm-GC.

1370
1371 *Dart*. Dart [13] compiles to Wasm-GC. It does not use `ref` and boxes scalars in a `struct`.
1372

1373 *Scheme*. The last addition to the small family of compiler targeting Wasm-GC is the Guile Scheme compiler. Scheme
1374 has many similar constraints as OCaml and Guile uses many similar solutions developed independently of Wasocaml.
1375 The compiler was presented to the Wasm-GC working group [26]. A more detailed explanation is published [25]. The
1376 implementation and its in-depth technical description are available [27].
1377
1378

1379 7.2 OCaml Web Compilers

1380 The history of OCaml compilers targeting web languages is quite crowded. This is surely thanks to the great pleasure
1381 that is writing compilers in OCaml.
1382
1383

1384 7.2.1 *Targeting JavaScript*. There are multiple OCaml compilers targeting JavaScript. Many approaches were experi-
1385 mented, the main two live ones are `js_of_ocaml` and `melange`. Naive compilation of OCaml to JavaScript is quite simple
1386 as it can almost be summarized as “type erasure”. There are some limitations in JavaScript that prevent that to be a
1387 complete compilation strategy. Moreover, getting to a proper compiler producing efficient and small JavaScript code
1388 is more complex.
1389

1390
1391 *Jsoo*. Jsoo [15] tries to be as close as possible to the native semantics. It compiles OCaml bytecode programs to
1392 JavaScript by decompiling it back to a simple untyped lambda calculus then performs many optimisation and minimi-
1393 sation passes.
1394

1395 *Melange*. Melange [8] tries to produce readable JavaScript, with a closer integration in the JavaScript module system.
1396 Melange starts from a modified version of the lambda IR which provides more type information. This allows the use
1397 of JavaScript features that match the uses of source features at the cost of some small semantical differences.
1398
1399

1400 7.2.2 Targeting Wasm.

1401
1402 *OCamlrun WebAssembly*. OCamlrun WebAssembly [7] compiles the OCaml runtime and interpreter to Wasm. They
1403 are both written in C so it is possible without any Wasm extension. In some cases, it allows to start executing code faster
1404

1405 than when compiling OCaml to Wasm. This is because there is less WebAssembly code to compile for the embedder.
1406 On the other hand, the execution is slower as it is interpreting OCaml bytecode inside Wasm instead of running a
1407 compiled version.
1408

1409 *WASICaml*. WASICaml [12] is quite similar to OCamlrun WebAssembly, the main difference being that WASICaml
1410 does not interpret the bytecode directly but translates it partially to Wasm. It leads to a faster execution.
1411

1412 These approaches suffer from the same problem as C programs running on Wasm: they cannot natively manipulate
1413 external values, such as DOM objects in the browser. Indeed, in the first versions of Wasm, the only types are scalars
1414 (**i32**, **i64**, **f32** and **f64**). There is no way to directly manipulate values from the embedder (e.g. JavaScript objects). It
1415 could be possible to identify objects with an integer and map them to their corresponding objects on the embedder side.
1416 This approach comes with the usual limitations of having two runtimes manipulating (indirectly) garbage-collected
1417 values: it is tedious and easily leads to memory leaks (cycles between the two runtimes cannot be collected).
1418
1419

1420 *Wasm_of_ocaml*. *Wasm_of_ocaml* [14] is a fork of *Js_of_ocaml*. It compiles the bytecode to Wasm-GC. It was de-
1421 veloped after *Wasocaml* and is based on the techniques we developed and that are described in the article. It does not
1422 benefit from the optimisations provided by Flambda. It is also quite restricted in its value representation choice and
1423 must use a representation based on arrays.
1424

1425 8 Conclusion

1426 The original goal of this implementation was to validate the Wasm-GC proposal and discuss its limits. This part was
1427 successful. No major changes were required for OCaml (mainly keeping the `i31ref` type and changes to the maximal
1428 length of sub-typing chains). And we are quite happy to claim that Wasm-GC is, to our opinion a very good design,
1429 that is a perfect compilation target for a garbage collected functional language like OCaml. And we think that the
1430 experience for most other garbage-collected languages will probably be similar.
1431

1432 As a side-effect, we now have an OCaml to Wasm-GC compiler. It is not yet usable because there are no user-side
1433 Wasm-GC runtime available. In order to test our compiler we are using V8 with various flags to enable experimental
1434 features such as the GC support. The design of the various extensions required by *Wasocaml* is stable and quite close
1435 to completion, but some details are still in flux. For this reason, we cannot expect them to be widely available soon. On
1436 the other hand, it means that our compiler will be ready when browsers start deploying new Wasm extensions.
1437

1438 Our future plans are to complete the Wasm implementations of OCaml externals, to implement the various FFI
1439 mechanisms, support effect handlers and to move *Wasocaml* to Flambda2. All of these will allow easy deployment of
1440 multi-language software on the browser while having good and predictable performances.
1441
1442
1443
1444

1445 References

- 1446 [1] Léo Andrès and Pierre Chambart. 2022. *Wasocaml*. <https://github.com/ocamlpro/wasocaml>
- 1447 [2] Léo Andrès and Pierre Chambart. 2023. OCaml on WasmGC. <https://github.com/WebAssembly/meetings/blob/main/gc/2023/GC-01-10.md>.
- 1448 [3] Andreas Haas, Andreas Rossberg, Derek L Schuff, Ben L Titzer, Michael Holman, Dan Gohman, Luke Wagner, Alon Zakai, and JF Bastien. 2017.
1449 Bringing the web up to speed with WebAssembly. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and*
1450 *Implementation*. 185–200.
- 1451 [4] Daniel Hillerström, Sam Lindley, Robert Atkey, and KC Sivaramakrishnan. 2017. Continuation passing style for effect handlers. (2017).
- 1452 [5] Donald E Knuth and Peter B Bendix. 1983. Simple word problems in universal algebras. *Automation of Reasoning: 2: Classical Papers on Computa-*
1453 *tional Logic 1967–1970* (1983), 342–376.
- 1454 [6] Chris Lattner. 2008. LLVM and Clang: Next generation compiler technology. In *The BSD conference*, Vol. 5. 1–20.
- 1455 [7] Sebastian Markbåge. 2017. *OCamlrun WebAssembly*. <https://github.com/sebmarkbage/ocamlrun-wasm>
- 1456 [8] Antonio Nuno Monteiro. 2022. *Melange*. <https://github.com/melange-re/melange>

- 1457 [9] Richard Musiol. 2018. WebAssembly architecture for Go. https://docs.google.com/document/d/131vj4DH6JFnb-blm_uRdaC0_
1458 [Nv3OUwjEY5qVCx4Cup4](https://docs.google.com/document/d/131vj4DH6JFnb-blm_uRdaC0_Nv3OUwjEY5qVCx4Cup4).
- 1459 [10] Gordon D Plotkin and Matija Pretnar. 2013. Handling algebraic effects. *Logical methods in computer science* 9 (2013).
- 1460 [11] Cheng Shao. 2022. WebAssembly backend merged into GHC. <https://www.tweag.io/blog/2022-11-22-wasm-backend-merged-in-ghc>.
- 1461 [12] Gerd Stolpmann. 2021. *WASICaml*. <https://github.com/remixlabs/wasicaml>
- 1462 [13] The Dart Project Authors. 2022. *dart2wasm*. <https://github.com/dart-lang/sdk/tree/main/pkg/dart2wasm>
- 1463 [14] Jérôme Vouillon. 2023. *Wasm_of_ocaml*. https://github.com/ocaml-wasm/wasm_of_ocaml
- 1464 [15] Jérôme Vouillon and Vincent Balat. 2014. From bytecode to JavaScript: the Js_of_ocaml compiler. *Software: Practice and Experience* 44, 8 (2014), 951–972.
- 1465 [16] WebAssembly Community Group participants. 2015. *Binaryen*. <https://github.com/WebAssembly/binaryen>
- 1466 [17] WebAssembly Community Group participants. 2017. *Exception Handling Proposal for WebAssembly*. <https://github.com/WebAssembly/exception-handling>
- 1467 [18] WebAssembly Community Group participants. 2017. *GC Proposal for WebAssembly*. <https://github.com/WebAssembly/gc>
- 1468 [19] WebAssembly Community Group participants. 2017. *Tail Call Proposal for WebAssembly*. <https://github.com/WebAssembly/tail-call>
- 1469 [20] WebAssembly Community Group participants. 2018. *Reference Types Proposal for WebAssembly*. <https://github.com/WebAssembly/reference-types>
- 1470 [21] WebAssembly Community Group participants. 2020. *Typed Function References Proposal for WebAssembly*. <https://github.com/WebAssembly/function-references>
- 1471 [22] WebAssembly Community Group participants. 2021. *JavaScript-Promise Integration Proposal for WebAssembly*. <https://github.com/WebAssembly/js-promise-integration>
- 1472 [23] WebAssembly Community Group participants. 2021. *Stack Switching Proposal for WebAssembly*. <https://github.com/WebAssembly/stack-switching>
- 1473 [24] WebAssembly Community Group participants. 2022. *Reference-Typed Strings Proposal for WebAssembly*. <https://github.com/WebAssembly/stringref>
- 1474 [25] Andy Wingo. 2023. a world to win: webassembly for the rest of us. <https://www.wingolog.org/archives/2023/03/20/a-world-to-win-webassembly-for-the-rest-of-us>.
- 1475 [26] Andy Wingo. 2023. Compiling Scheme to WasmGC. <https://github.com/WebAssembly/meetings/blob/main/gc/2023/GC-04-18.md>.
- 1476 [27] Andy Wingo. 2023. *Hoot*. <https://gitlab.com/spritely/guile-hoot/-/blob/main/design/ABI.md>
- 1477
- 1478
- 1479
- 1480
- 1481
- 1482
- 1483
- 1484
- 1485
- 1486
- 1487
- 1488
- 1489
- 1490
- 1491
- 1492
- 1493
- 1494
- 1495
- 1496
- 1497
- 1498
- 1499
- 1500
- 1501
- 1502
- 1503
- 1504
- 1505
- 1506
- 1507
- 1508