

HPVM2Wasm workflow

WASM modules
(host code + data
flow code)

→ HPVM binaries

Wasm modules
(cpu host code) +
wasm modules
webgpu modules



HPVM
graph
structure
(.json)

+

HPVM2Web
runtime
(emscripten?)

compiler infrastructure for WASM

- simple C++ API header
- can be used from JS
- accepts WASM-like input but also control flow graph
- internal IR designed for parallel code gen
- used to build compilers
- provides toolchain utilities that can...
- parse, optimize, and re-emit WASM
- interpret WASM
- integrate with emscripten for C++ → WASM
- binaries IR
- very close to WASM, but a tree not a stack machine
- S-expr based

JIT compilation
in browser

RUNTIME

emcc - rt.js included
in all the example dirs

- built on top of emscripten runtime for web assembly
 - JS runtime environment
- implements subset of execution semantics of HPVM data flow graph
 - nodes are when ready
 - hierarchical DAG execution
 - dynamic node replication
- each HPVM leaf node is a separate WASM module
 - node JIT compiled when ready
 - custom memory tracker and management for copying memory
- nodes executed asynchronously using JS promises
- automatic compilation support for synchronous C/C++ code with async JS
 - emscripten asyncify compilation flow
 - overload HPVM host code API to invoke JS



I/O

- uses SDL for browser environment
- applications that use SDL require no change

File systems

- browser prevents direct host file access
- emscripten provides own libc and libexec
- creates own virtual file system

Browser main loop

- each event has a turn and then returns control (cooperative multitasking)

Execution lifecycle

- preloading phase:

- run phase

Memory Representation

- pointers at offsets

using
HPVM binaries
compiler GPU
code gen

don't understand
this

allows sync C++ code to interact w/ async JS

- transforms compiled code into code that can be paused and resumed
- handles the pausing + resuming
- implemented as a binaries pass
- alternatives: coroutines (not supported by WASM) and threading (not supported by all browsers)
- sleep actually returns control (for example)

Graph Description

- dataflow graph encoded in json
- each leaf node specifies Wasm module
- json ported by runtime, which manages heterogeneous compilation + memory

BN the user writes but ideally generated
HPVM could generate easily

Host code API

- defined entirely in JS and imported into host code module with Wasm

defined in hpvm2wasmRT.js

Wasm to WGSL

webGPU shading language

- implemented as Binaryen pass
 - Aaron implemented Binaryen pass
 - operates on the s-expr format of Binaryen IR
 - takes wasm to a CFG
 - identifies patterns to build block structure for webGPU (big task #1)
 - webGPU doesn't support pointer arith, rewrite pointer access (big task #2)
 - this is compiled to binaryen.js
- translates a single specified function to a webGPU kernel

Aaron implemented pass that operates on s-expr

- takes wasm → CFG
- tries to identify patterns
- to blocks for webGPU
- block structures?
- pointer access

in other words...

translation process:

1. construct control flow graph
2. compute dominance (?)
3. identify back edges + loops
4. identify if-then-else
5. generate code for each basic block
 - array access reconstruction (?)
6. glue basic blocks together w/ appropriate control flow

happens w/in
binaryen pass

WebGPU runtime

- JS code
- creates buffers for each array argument and buffer of scalar values
- produces shader module
- binds buffer into the shader (copies memory over?)
- executes the shader
- copies data back from GPU

nodes can be marked for gpu - wasm runtime considers this marking + may send node here (in execute_node())

Future work

- additional accelerator backends
- integration with browsers?
- compilation from hetero c++ to hpvm2wasm
- process virtualization (same physical machine but isolated envs)
- dynamically targeting nodes to devices (at runtime?)
- other optimizations before code gen

have gpu + cpu support but not for other accel. (like ML accel)

Translating hetero-c++ to hpvm2web

Big Action Items:

- 1) generate .json from HPVM
- 2) kernels in separate files... is this necessary?
- 3) lower HPVM mnemonics

Project for the rest of semester:

- hetero-c++ to web assembly for CPU only
- involves creating the wasm modules and integrating the JS runtime
- probably need to do some work in the runtime

What does the hpvm2web matmul look like?

- 1) matmul.json
 - describes the graph (should be easy to generate from HPVM)