



Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

Life of a Taichi Kernel

A trip through Taichi's internal design and implementation

Yuanming Hu

August 24, 2020



Before we start: interacting with this file

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

This PDF file has links that point you to more details:

- C++ source code, e.g., ▶[program/program.h](#)
- Python source code, e.g., ▶[lang/kernel.py](#)
- Documentation, e.g., ▶[doc:Hello, world!](#)



Overview

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

This talk (or design doc) briefly covers some internal design decisions of Taichi, for developers and users who want to dig deeper into Taichi. Topics include

- How does Taichi work?
- Why do we end up the current design decisions?
- What do I need to know to improve Taichi?

Setting up Taichi locally

- Installation via pip: `python3 -m pip install taichi`
- Building from source (or Docker) for development: [▶doc:Developer installation](#)

Pick one installation only

Avoid having both pip-installed Taichi and the version build from source.



Table of Contents

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

- 1 Introduction
- 2 Python frontend
- 3 Intermediate representation
 - Computation IR
 - Structural node IR
- 4 Sparse programming
- 5 Differentiable programming
- 6 Backends
 - Runtime system
- 7 Summary



What is Taichi?

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

High-performance domain-specific language (DSL) embedded in **Python**, for **computer graphics** applications

- **Productivity** and **portability**: easy to learn, to write, and to share
- **Performance**: data-oriented, parallel, megakernels
- **Spatially sparse** programming: save computation and storage on empty regions
- **Decouple** data structures from computation
- **Differentiable** programming support



Taichi v.s. deep learning frameworks

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

Why is Taichi different from TensorFlow, PyTorch, NumPy, JAX, ... ?

Quick answer: Taichi uniquely supports **megakernels** and **spatial sparsity**.

Longer answer: Those systems serve their own application domains (e.g., convolutional neural networks) very well, but their design decisions surrounding **immutable, dense tensors** (e.g., feature maps) with **simple, regular operators** (e.g., element-wise add and 2D convolutions) do not serve well more irregular computational patterns, including

- Computer graphics, including physical simulation and rendering
- Irregular neural network layers (e.g., gathering/scattering) that are emerging
- General differentiable programming cases

Without Taichi people tend to manually write CUDA or abuse deep learning programming interfaces. Taichi offers performance, productivity, and portability in those cases.



Hello, world! (Julia set, $z \leftarrow z^2 + c$)

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

```
import taichi as ti
ti.init(arch=ti.gpu)
n = 320
pixels = ti.field(dtype=float, shape=(n * 2, n))

@ti.func
def complex_sqr(z):
    return ti.Vector([z[0]**2 - z[1]**2, z[1] * z[0] * 2])

@ti.kernel
def paint(t: float):
    for i, j in pixels: # Parallized over all pixels
        c = ti.Vector([-0.8, ti.cos(t) * 0.2])
        z = ti.Vector([i / n - 1, j / n - 0.5]) * 2
        iterations = 0
        while z.norm() < 20 and iterations < 50:
            z = complex_sqr(z) + c
            iterations += 1
        pixels[i, j] = 1 - iterations * 0.02

gui = ti.GUI("Julia Set", res=(n * 2, n))

for i in range(1000000):
    paint(i * 0.03)
    gui.set_image(pixels)
    gui.show()
```

More details: ▶[doc:Hello, world!](#) Run it: `ti example fractal`



Table of Contents

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

- 1 Introduction
- 2 Python frontend
- 3 Intermediate representation
 - Computation IR
 - Structural node IR
- 4 Sparse programming
- 5 Differentiable programming
- 6 Backends
 - Runtime system
- 7 Summary



Taichi's frontend

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

Taichi's old C++ frontend

Taichi used to be embedded in C++14. However, that solution is mostly abandoned, because

- ① C++ itself is too complex for most users to learn, not to say a DSL embedded in C++. E.g., `math/svd.h` 😞
- ② Getting C++ AST is almost impossible. We had to heavily use templates/macros tricks, which harm readability. 😞

Taichi's new Python frontend

Now the whole Taichi system is deeply embedded in Python.

- ① Python is easy to learn and widely adopted. 😊
- ② Python allows flexible AST inspection and manipulation. 😊

Both of these allow us to invent a new high-performance language out of Python.



Other goodies of the Python frontend

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

- Easy to run. No ahead-of-time compilation is needed.
- Reuse and interact with existing python infrastructure:
 - ① IDEs such as PyCharm.
 - ② Package manager (pip)
 - ③ Existing packages such as `matplotlib` and `numpy`
- The built-in AST manipulation tools (`import ast`) in python allow us to do magical things, as long as the kernel body can be parsed by the Python parser.

Kernels marked with `@ti.kernel` will be compiled into a Taichi AST and then to parallel kernels on CPU/GPUs.



Generating a Taichi AST from Python AST

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

(i) Taichi kernel (Python) \rightarrow Taichi AST generator (Python)

Taichi has a series of Python AST transformers `►lang/transformer.py` that transforms a Taichi kernel (in Python) into another Python script, which is a Taichi AST generator.

(ii) Taichi AST generator (Python) \rightarrow Taichi AST (C++)

The Taichi AST generator is a Python script that calls AST builder functions in C++ (exported via pybind11) when executed. The result of this step is a Taichi AST.

Confused? Let's take a look at an example 😊



Generating a Taichi AST (example)

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

Example

```
import taichi as ti

ti.init(print_preprocessed=True)

@ti.kernel
def foo():
    for i in range(10):
        if i == 2:
            print(i)

foo()
```

Inspecting Taichi AST transforms

Set `print_preprocessed=True` in `ti.init` to make Taichi print out processed AST.



Generating a Taichi AST (example)

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

```
# (1) Input kernel (Python function)
@ti.kernel
def foo():
    for i in range(10):
        if i % 2 == 0:
            print(i)
```

```
# (2) Transformed AST generator (another Python function)
def foo():
    import taichi as ti
    if 1:
        i = ti.Expr(ti.core.make_id_expr(''))
        __begin = ti.Expr(0)
        __end = ti.Expr(10)
        __begin = ti.cast(__begin, ti.i32)
        __end = ti.cast(__end, ti.i32)
        ti.core.begin_frontend_range_for(i.ptr,
            __begin.ptr, __end.ptr)
        if 1:
            __cond = ti.chain_compare([i, 2], ['Eq'])
            ti.core.begin_frontend_if(ti.Expr(__cond).ptr)
            ti.core.begin_frontend_if_true()
            ti.ti_print(i)
            ti.core.pop_scope()
            ti.core.begin_frontend_if_false()
            ti.core.pop_scope()
        ti.core.end_frontend_range_for()
    del i # Note: Taichi has lexical scoping!
```

```
# (3) Generated Taichi AST
kernel {
    $0 : for @tmp0 in range(
        (cast_value<int32> 0),
        (cast_value<int32> 10))
        block_dim=adaptive {
        $1 : if (1 & (@tmp0 == 2)) {
            $2 = eval @tmp0
            print %2, "\n"
        } else {
        }
        }
    }
}
```

Note

Taichi (frontend) AST \neq
Taichi (middle-end) IR



Table of Contents

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

- 1 Introduction
- 2 Python frontend
- 3 Intermediate representation**
 - Computation IR
 - Structural node IR
- 4 Sparse programming
- 5 Differentiable programming
- 6 Backends**
 - Runtime system
- 7 Summary



Taichi intermediate representation (IR): Two components

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR

Structural node IR

Sparse

programming

Differentiable
programming

Backends

Runtime system

Summary

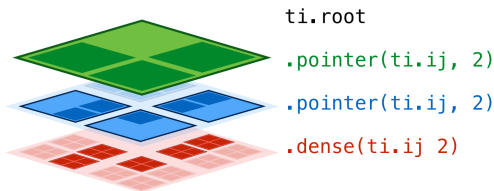
1. Computation IR

- Static-single assignment
- Hierarchical (instead of CFG+BB)
- Differentiable
- Statically and strongly typed

```
$0 = offloaded range_for(0, 10)
body {
  <i32 x1> $1 = loop $0 index 0
  <i32 x1> $2 = const [3]
  <i32 x1> $3 = cmp_eq $1 $2
  <i32 x1> $4 = const [1]
  <i32 x1> $5 = bit_and $3 $4
  $6 : if $5 {
    print $1, "\n"
  }
}
```

2. Structural Node (SNode) IR

- Describes data organization
- Tree-structured
- Spatial sparsity
- Favors powers of two



IR design goals: a) enable domain-specific optimizations for sparse computation.
b) decouple data layout from computation (more on this later).



Progressive compilation

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

- 1 Taichi compiler progressively compiles a frontend AST into executable CPU/GPU kernels.
- 2 During this process, 30+ compilation passes are applied to the Taichi AST/IR.
- 3 The only big step is from Taichi AST to Taichi IR (AST lowering ▶ `transforms/lower_ast.cpp`), where two different sets of statements (instructions) are used.
- 4 Each other pass makes a small step towards hardware-friendly code only. Almost the same sets of statements are used for input/output Taichi IR.
- 5 There is **no** clear separation of IR statements (e.g. high-level, mid-level, low-level IR) after AST lowering.



10 key compilation passes ▶ [transforms](#)

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

Taichi's IR transformation passes gradually converts an input Taichi AST to parallel executable kernels ▶ [transforms/compile_to_offloads.cpp](#).

① Taichi frontend AST to Taichi IR

- ① Lower Taichi AST to SSA ▶ [transforms/lower_ast.cpp](#)
- ② Type checking ▶ [transforms/type_check.cpp](#)

② Taichi IR to offloaded tasks

- ① (Optional) Automatic differentiation ▶ [transforms/auto_diff.cpp](#)
- ② (Optional) Insert bound checks ▶ [transforms/check_out_of_bound.cpp](#)
- ③ Flag and weaken access ▶ [transforms/flag_access.cpp](#)
- ④ Automatic parallelization ▶ [transforms/offload.cpp](#)

③ Offloaded tasks to executable

- ① Demote dense struct-fors to range-fors
▶ [transforms/demote_dense_struct_fors.cpp](#)
- ② Create thread local storage ▶ [transforms/make_thread_local.cpp](#)
- ③ Create block local storage ▶ [transforms/make_block_local.cpp](#)
- ④ Lower access ▶ [transforms/lower_access.cpp](#)

(Optimization passes are skipped to save space)



Statements

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

Taichi IR has ~ 70 statements.

Typical statements include:

- Arithmetic ▶ [ir/ir.h](#)
 - ① UnaryOpStmt. Operators ▶ [inc/unary_op.inc.h](#) (sqrt, sin, ...)
 - ② BinaryOpStmt. Operators ▶ [inc/binary_op.inc.h](#) (add, mul, ...)
 - ③ ...
- Memory access: `Global[Ptr/Load/Store]Stmt`
- AutoDiff Stack operations (more on this later):
`Stack[Alloca/LoadTop/LoadTopAdj/Pop/Push/AccAdjoint]Stmt`
- SNode Micro Ops (after the `lower_access` pass)
- ...



General-purpose optimization passes

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

Why optimize Taichi IR? Can't we leave the job to LLVM?

Quick answer:

- Higher optimization quality (next slide)
- Fewer instructions to LLVM: faster JIT

General-purpose optimization passes

- Control-flow based optimizations: CSE, DIE, ...
▶ [transforms/cfg_optimization.cpp](#)
- Constant folding ▶ [transforms/constant_fold.cpp](#) (more details later)
- Algebraic simplification ▶ [transforms/alg_simp.cpp](#)
- ...



Why not fully trust backend compilers?

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

Long story short: backend compilers (e.g., LLVM) do not have as much information as the Taichi compiler do.

For example, Taichi IR and optimizer

- ① ...can do index analysis
- ② ...have a tailored instruction granularity
- ③ ...have strict data access semantics

Therefore traditional compilers often fail to do certain important optimizations.

Benchmarks show Taichi IR optimization passes (mostly on data accesses) make programs $3.02\times$ faster¹, even if we use `-O3` on backend compilers.

¹Y. Hu et al. (2019). “Taichi: a language for high-performance computation on spatially sparse data structures”. In: *ACM Transactions on Graphics (TOG)* 38.6, pp. 1–16.



The IR granularity spectrum

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

$x[i, j]$

access1(i,j)
access2(i,j)

```
$4 = [$4][root]::lookup(root, $3) coord = ($2)
  activate = false
$5 = get_child [$4->$3] $4
$6 = bit_extract($2 + 0, 7-14)
$7 = linearized(ind [$6], stride (128))
$8 = [$3][dense]::lookup($5, $7) coord = ($2)
  = false
$9 = get_child [$3->$2] $8
$10 = bit_extract($2 + 0, 0-7)
$11 = linearized(ind [$10], stride (128))
$12 = [$2][dense]::lookup($9, $11) coord = ($2)
  activate = false
```

```
%63 = lshr i32 %62, 0
%64 = and i32 %63, 255
%65 = add i32 %63, 0
%66 = lshr i32 %65, 0
%67 = and i32 %66, 255
%68 = add i32 0, %64
%69 = mul i32 %68, 256
%70 = add i32 %69, %67
%71 = bitcast %struct.DenseMeta* %5 to
call void @StructMeta_set_snode_id(%5,
call void @StructMeta_set_element_size
call void @StructMeta_set_max_num_ele
call void @StructMeta_set_lookup_elem
element)
call void @StructMeta_set_is_active(%5
call void @StructMeta_set_get_num_ele
lements)
call void @StructMeta_set_from_parent.
```

```
movl $0, %eax
addl %eax, %ebx
popl %eax
looptop:
imul %edx
andl $0xFF, %eax
cmpl $100, %eax
jb looptop
leal 4(%esp), %ebp
movl %esi, %edi
subl $8, %edi
shrl %cl, %ebx
movw %bx, -2(%ebp)
movl $0, %eax
addl %eax, %ebx
popl %eax
looptop:
imul %edx
andl $0xFF, %eax
cmpl $100, %eax
jb looptop
leal 4(%esp), %ebp
movl %esi, %edi
subl $8, %edi
shrl %cl, %ebx
movw %bx, -2(%ebp)
```

End2end access

Level-wise Access

Taichi IR

LLVM IR

Machine code

Coarser

Finer



The IR granularity trade-offs

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

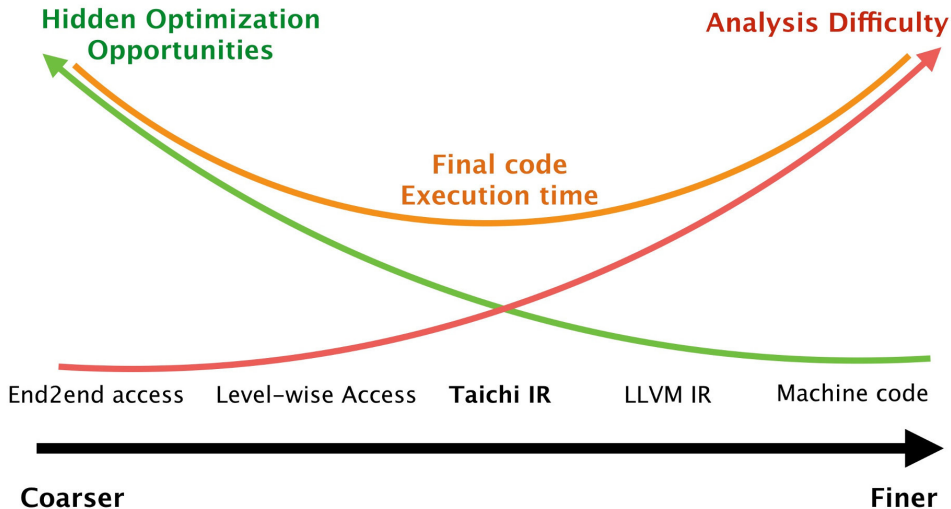
Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary





Data access semantics in Taichi

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

Seemingly trivial assumptions can enable effective compiler optimization and make programmers' life easier. E.g.,

- ① No pointer aliasing: $a[x, y]$ and $b[i, j]$ never overlaps unless a and b are the same field ▶ [analysis/alias_analysis.cpp](#)
- ② All memory accesses are done through `field[indices]` syntax
 - Pointers that can flexibly point to anything arguably makes optimization harder
- ③ The only way data structures get modified, is through write accesses of form `field[indices]`
- ④ Read access **does not** modify anything
 - No memory allocation
 - No exception if element does not exist



Example general-purpose opt pass: Constant folding

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

Taichi's constant folding pass uses JIT evaluation on the device

► [transforms/constant_fold.cpp](#)

- 1 Generate a small Taichi kernel that only does the simple unary/binary operation. E.g, `f32 + f32`.
- 2 Evaluate the small kernel on operands
- 3 Replace the arithmetic statement with a constant statement with value being the result.

Evaluating on the device for correctness

Certain operations (e.g., `sin(x)`) on host evaluation (e.g., x64) may generate a different result from device (e.g., OpenGL) evaluation.

For correctness, do constant folding evaluation on devices. (... and assume no `fastmath`.)



Constant folding (example)

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

Madhava's series:

$$\pi = \sqrt{12} \sum_{k=0}^{\infty} \frac{(-3)^{-k}}{2k+1} = \sqrt{12} \sum_{k=0}^{\infty} \frac{(-\frac{1}{3})^k}{2k+1} = \sqrt{12} \left(1 - \frac{1}{3 \cdot 3} + \frac{1}{5 \cdot 3^2} - \frac{1}{7 \cdot 3^3} + \dots \right)$$

```
import taichi as ti
ti.init(print_ir=True)

@ti.kernel
def calc_pi() -> ti.f32:
    s = 0.0
    c = 1.0
    for i in ti.static(range(10)):
        s += c / (i * 2 + 1)
        c *= -1/3
    return s * ti.sqrt(12.0)

print(calc_pi())
```



Constant folding (example)

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

Madhava's series:

$$\pi = \sqrt{12} \sum_{k=0}^{\infty} \frac{(-3)^{-k}}{2k+1} = \sqrt{12} \sum_{k=0}^{\infty} \frac{(-\frac{1}{3})^k}{2k+1} = \sqrt{12} \left(1 - \frac{1}{3 \cdot 3} + \frac{1}{5 \cdot 3^2} - \frac{1}{7 \cdot 3^3} + \dots \right)$$

Generated Taichi IR:

```
kernel {  
  $0 = offloaded {  
    <f32 x1> $1 = const [3.1415904] # All computation folded  
    <f32 x1> $2 : kernel return $1  
  }  
}
```



Offloading (automatic parallelization)

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

The offload pass ▶[transforms/offload.cpp](#) decomposes kernels into **offloaded statements** (`OffloadedStmt`). Each `OffloadedStmt` have one of the following types:

- ① `serial`: Simple serial code
- ② `range_for`: Parallel range-for
- ③ `clear_list`: Clear the list of nodes for `struct_for`
- ④ `listgen`: Generate list of nodes for `struct_for`
- ⑤ `struct_for`: Parallel struct-for
- ⑥ `gc`: Garbage collection



Offloading (automatic parallelization transform)

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

Challenge: outside variables are visible to loop bodies.

Solution: create global temporary variables (`global tmp var`)

```
import taichi as ti

ti.init(
    print_ir=True,
    use_thread_local=False
)

@ti.kernel
def bar():
    a = 0
    for i in range(10):
        a += i
    print(a)

bar()
```

Note: `a` is accessible in the loop body (parallel), but its declaration (serial)/printing (serial) is outside.

```
kernel {
    $0 = offloaded serial
    body {
        <i32*x1> $1 = global tmp var (offset = 0 B)
        <i32 x1> $2 = const [0]
        <i32*x1> $3 : global store [$1 <- $2]
    }
    $4 = offloaded serial range_for(0, 10) grid_dim=0
        block_dim=32
    body {
        <i32 x1> $5 = loop $4 index 0
        <i32*x1> $6 = global tmp var (offset = 0 B)
        <i32 x1> $7 = atomic add($6, $5)
    }
    $8 = offloaded serial
    body {
        <i32*x1> $9 = global tmp var (offset = 0 B)
        <i32 x1> $10 = global load $9
        print $10, "\n"
    }
}
```



Block local storage (BLS)

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

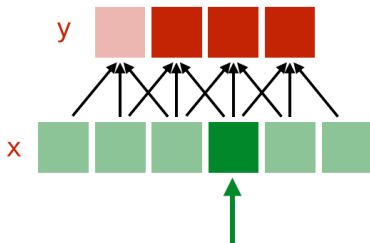
Differentiable
programming

Backends
Runtime system

Summary

► `transforms/make_block_local.cpp`

```
for i in x:  
     $y[i] = x[i - 1] - 2 * x[i] + x[i + 1]$ 
```



3x global memory accesses for a single element

**To reduce global memory accesses
(using block local storage):**

Each block first fetches all the needed **x**'s to a **block local** buffer. When evaluating **y**, read from the block local buffer instead of from global memory.



Block local storage (BLS): Workflow on CUDA

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

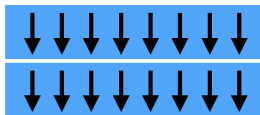
Differentiable
programming

Backends

Runtime system

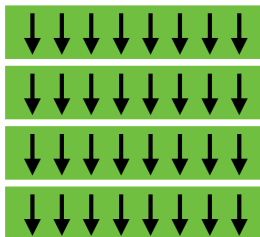
Summary

Prologue
fetch/zero-fill BLS buffer
(**global** to **shared**)



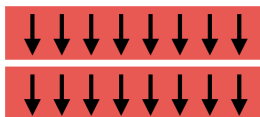
—syncthreads

Body



—syncthreads

Epilogue (optional)
Write-back BLS buffer
(**shared** to **global**)





Block local storage (BLS): Productivity

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR

Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

CUDA

```
__global__ void laplace_shared(float *a, float *b) {
    unsigned int i = blockIdx.x * blockDim.x + threadIdx.x;
    unsigned int j = blockIdx.y * blockDim.y + threadIdx.y;

    unsigned int tid = blockDim.y * threadIdx.x + threadIdx.y;

    __shared__ float pad[bs + 2][bs + 2];
    auto pad_size = (bs + 2) * (bs + 2);

    if (bs <= i && i < N - bs && bs <= j && j < N - bs) {
        while (tid < pad_size) {
            int si = tid / (bs + 2);
            int sj = tid % (bs + 2);
            int gi = si - 1 + blockIdx.x * blockDim.x;
            int gj = sj - 1 + blockIdx.y * blockDim.y;
            pad[si][sj] = a[gi * N + gj];
            tid += blockDim.x * blockDim.y;
        }
    }

    __syncthreads();

    if (bs <= i && i < N - bs && bs <= j && j < N - bs) {
        auto ret = -4 * pad[threadIdx.x + 1][threadIdx.y + 1] +
            pad[threadIdx.x + 2][threadIdx.y + 1] +
            pad[threadIdx.x][threadIdx.y + 1] +
            pad[threadIdx.x + 1][threadIdx.y + 2] +
            pad[threadIdx.x + 1][threadIdx.y];

        b[i * N + j] = ret;
    }
}
```

Taichi

```
@ti.kernel
def laplace():
    ti.cache_shared(a)
    for i, j in a:
        b[i, j] = -4 * a[i, j] - a[i - 1][j] - \
            a[i][j - 1] - a[i + 1][j] - a[i][j + 1]
```



Block local storage (BLS): Productivity

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR

Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

CUDA

```
1 // CUDA kernel for particle state update and scatter to grid (P2G)
2 // This kernel is generated by Taichi Python frontend.
3 // It is a block local storage (BLS) kernel.
4 // The kernel is written in C++ and compiled to CUDA.
5 // The kernel is a template function.
6 // The kernel is a block local storage (BLS) kernel.
7 // The kernel is a block local storage (BLS) kernel.
8 // The kernel is a block local storage (BLS) kernel.
9 // The kernel is a block local storage (BLS) kernel.
10 // The kernel is a block local storage (BLS) kernel.
11 // The kernel is a block local storage (BLS) kernel.
12 // The kernel is a block local storage (BLS) kernel.
13 // The kernel is a block local storage (BLS) kernel.
14 // The kernel is a block local storage (BLS) kernel.
15 // The kernel is a block local storage (BLS) kernel.
16 // The kernel is a block local storage (BLS) kernel.
17 // The kernel is a block local storage (BLS) kernel.
18 // The kernel is a block local storage (BLS) kernel.
19 // The kernel is a block local storage (BLS) kernel.
20 // The kernel is a block local storage (BLS) kernel.
21 // The kernel is a block local storage (BLS) kernel.
22 // The kernel is a block local storage (BLS) kernel.
23 // The kernel is a block local storage (BLS) kernel.
24 // The kernel is a block local storage (BLS) kernel.
25 // The kernel is a block local storage (BLS) kernel.
26 // The kernel is a block local storage (BLS) kernel.
27 // The kernel is a block local storage (BLS) kernel.
28 // The kernel is a block local storage (BLS) kernel.
29 // The kernel is a block local storage (BLS) kernel.
30 // The kernel is a block local storage (BLS) kernel.
31 // The kernel is a block local storage (BLS) kernel.
32 // The kernel is a block local storage (BLS) kernel.
33 // The kernel is a block local storage (BLS) kernel.
34 // The kernel is a block local storage (BLS) kernel.
35 // The kernel is a block local storage (BLS) kernel.
36 // The kernel is a block local storage (BLS) kernel.
37 // The kernel is a block local storage (BLS) kernel.
38 // The kernel is a block local storage (BLS) kernel.
39 // The kernel is a block local storage (BLS) kernel.
40 // The kernel is a block local storage (BLS) kernel.
41 // The kernel is a block local storage (BLS) kernel.
42 // The kernel is a block local storage (BLS) kernel.
43 // The kernel is a block local storage (BLS) kernel.
44 // The kernel is a block local storage (BLS) kernel.
45 // The kernel is a block local storage (BLS) kernel.
46 // The kernel is a block local storage (BLS) kernel.
47 // The kernel is a block local storage (BLS) kernel.
48 // The kernel is a block local storage (BLS) kernel.
49 // The kernel is a block local storage (BLS) kernel.
50 // The kernel is a block local storage (BLS) kernel.
51 // The kernel is a block local storage (BLS) kernel.
52 // The kernel is a block local storage (BLS) kernel.
53 // The kernel is a block local storage (BLS) kernel.
54 // The kernel is a block local storage (BLS) kernel.
55 // The kernel is a block local storage (BLS) kernel.
56 // The kernel is a block local storage (BLS) kernel.
57 // The kernel is a block local storage (BLS) kernel.
58 // The kernel is a block local storage (BLS) kernel.
59 // The kernel is a block local storage (BLS) kernel.
60 // The kernel is a block local storage (BLS) kernel.
61 // The kernel is a block local storage (BLS) kernel.
62 // The kernel is a block local storage (BLS) kernel.
63 // The kernel is a block local storage (BLS) kernel.
64 // The kernel is a block local storage (BLS) kernel.
65 // The kernel is a block local storage (BLS) kernel.
66 // The kernel is a block local storage (BLS) kernel.
67 // The kernel is a block local storage (BLS) kernel.
68 // The kernel is a block local storage (BLS) kernel.
69 // The kernel is a block local storage (BLS) kernel.
70 // The kernel is a block local storage (BLS) kernel.
71 // The kernel is a block local storage (BLS) kernel.
72 // The kernel is a block local storage (BLS) kernel.
73 // The kernel is a block local storage (BLS) kernel.
74 // The kernel is a block local storage (BLS) kernel.
75 // The kernel is a block local storage (BLS) kernel.
76 // The kernel is a block local storage (BLS) kernel.
77 // The kernel is a block local storage (BLS) kernel.
78 // The kernel is a block local storage (BLS) kernel.
79 // The kernel is a block local storage (BLS) kernel.
80 // The kernel is a block local storage (BLS) kernel.
81 // The kernel is a block local storage (BLS) kernel.
82 // The kernel is a block local storage (BLS) kernel.
83 // The kernel is a block local storage (BLS) kernel.
84 // The kernel is a block local storage (BLS) kernel.
85 // The kernel is a block local storage (BLS) kernel.
86 // The kernel is a block local storage (BLS) kernel.
87 // The kernel is a block local storage (BLS) kernel.
88 // The kernel is a block local storage (BLS) kernel.
89 // The kernel is a block local storage (BLS) kernel.
90 // The kernel is a block local storage (BLS) kernel.
91 // The kernel is a block local storage (BLS) kernel.
92 // The kernel is a block local storage (BLS) kernel.
93 // The kernel is a block local storage (BLS) kernel.
94 // The kernel is a block local storage (BLS) kernel.
95 // The kernel is a block local storage (BLS) kernel.
96 // The kernel is a block local storage (BLS) kernel.
97 // The kernel is a block local storage (BLS) kernel.
98 // The kernel is a block local storage (BLS) kernel.
99 // The kernel is a block local storage (BLS) kernel.
100 // The kernel is a block local storage (BLS) kernel.
```

Taichi

```
ti.cache_shared(grid_m, grid_v)
for I in ti.grouned(pid): # Particle state update and scatter to grid (P2G)
    p = pid[I]
    base = ((x[p] + Offset) * inv_dx - 0.5).cast(int)
    for d in ti.static(range(3)):
        base[d] = ti.assume_in_range(base[d], I[d], 0, 1)
    fx = (x[p] + Offset) * inv_dx - base.cast(float)
    w = [0.5 * ti.sqr(1.5 - fx), 0.75 - ti.sqr(fx - 1), 0.5 * ti.sqr(fx - 0.5)]
    stress = ...
    stress = (-dt * p_vol * 4 * inv_dx * inv_dx) * stress
    affine = stress + p_mass * C[p]
    for i, j, k in ti.static(ti.ndrange(3, 3, 3)): # Loop over 3x3x3 grid
        offset = ti.Vector([i, j, k])
        dpos = (offset.cast(float) - fx) * dx
        weight = w[i][0] * w[j][1] * w[k][2]
        grid_v[base + offset] += weight * (p_mass * v[p] + affine @ dpos)
        grid_m[base + offset] += weight * p_mass
```




Data types

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

Taichi is statically and strongly and typed. Supported types include

- Signed integers: `ti.i8/i16/i32/i64`
- Unsigned integers: `ti.u8/u16/u32/u64`
- Float-point numbers: `ti.f32/f64`

`ti.i32` and `ti.f32` are the most commonly used types in Taichi. Boolean values are represented by `ti.i32` for now.

Data type compatibility

The CPU and CUDA backends support all data types. Other backend may miss certain data type support due to backend API constraints. See the documentation for more details.



Fields

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

Taichi is a *data-oriented* programming language where **fields** are first-class citizens.

- Fields are essentially multi-dimensional arrays
- An element of a field can be either a scalar (`ti.field`), a vector (`ti.Vector.field`), or a matrix (`ti.Matrix.field`)
- Field elements are *always* accessed via the `a[i, j, k]` syntax. (No pointers.)
- Access out-of-bound is undefined behavior in non-debug mode
- (*Advanced*) Fields can be spatially sparse



Manipulating data layouts

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

- **Data layout** refers to the dimensionality, shape, memory layout, `and` sparsity structure of fields.
- A carefully designed data layout can significantly improve cache/TLB-hit rates and cacheline utilization.
- Taichi decouples algorithms from data layouts, and the Taichi compiler automatically optimizes data accesses on a specific data layout.
- These Taichi features allow programmers to quickly experiment with different data layouts and figure out the most efficient one on a specific task and computer architecture.
- In Taichi, the layout is defined in a recursive manner.

(From now on we focus on scalar fields `ti.field()`.)



Data layout in three languages

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR

Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

Taichi (Python)

```
x = ti.field(dtype=ti.i32);  
ti.root.dense(ti.i, 16).place(x)
```

Equivalent C++

```
int x[16];
```

Natural language

- `ti.root` is the root data structure.
- Dot (`A.B`) means “Each cell of `A` has `B`”
- `dense(ti.i, 16)` means “a `dense` container with 16 cells along the `ti.i` axis”
- `place` means “... `field` ...”

“Each cell of `root` has a `dense` container with 16 cells along the `ti.i` axis.
Each cell of `dense` has `field x`.”



More data layouts

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

Structure of arrays (SOA):

Taichi (Python)

```
ti.root.dense(ti.i, 16).place(x)
ti.root.dense(ti.i, 16).place(y)
```

Array of structures (AOS):

Taichi (Python)

```
ti.root.dense(ti.i, 16).place(x, y)
# or equivalently
point = ti.root.dense(ti.i, 16)
point.place(x)
point.place(y)
```

Equivalent C++

```
int x[16];
int y[16];
```

Equivalent C++

```
struct Point {int x, y}
Point points[16];
```



Nesting: Array of structures of arrays (AOSOA) example

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

Taichi (Python)

```
chunk = ti.root.dense(ti.i, 2)

chunk.dense(ti.i, 8).place(x)
chunk.dense(ti.i, 8).place(y)
```

Equivalent C++

```
struct PointChunk {
    int x[8]; int y[8];
};
PointChunk points[2];
```

Discussions

In Taichi, you can always access x and y using syntax like $x[5]$ or $y[7]$.
In C++ the access is data-layout dependent.

This is just for dense data. For sparse data it is even trickier in C++.



Higher dimensionality

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

Taichi (Python)

```
ti.root.dense(ti.ij, (4, 8)).place(x)
```

Taichi (Python)

```
block = ti.root.dense(ti.ij, 2)
block.dense(ti.ij, 4).place(x)
# or simply
ti.root.dense(ti.ij, 2)
    .dense(ti.ij, 4).place(x)
```

Equivalent C++

```
int x[4][8];
```

Equivalent C++

```
struct Block {
    int x[4][4];
};
Block blocks[2][2];
```



Structural Nodes (SNodes) ▶ [ir/snode.h](#)

Life of a Taichi
Kernel

Yuanming Hu

[Introduction](#)

[Python frontend](#)

[Intermediate
representation](#)

Computation IR
Structural node IR

[Sparse
programming](#)

[Differentiable
programming](#)

[Backends](#)

Runtime system

[Summary](#)

Currently supported SNode types ▶ [inc/snodes.inc.h](#):

- ① **root**: The root of the data structure.
- ② **dense**: A fixed-length contiguous array.
- ③ **bitmasked**: similar to **dense**, but it also uses a mask to maintain sparsity information, one bit per child.
- ④ **pointer**: Store pointers instead of the whole structure to save memory and maintain sparsity.
- ⑤ **dynamic**: Variable-length array, with a predefined maximum length. It serves the role of `std::vector` in C++ or **list** in Python.



Access lowering [▶transforms/lower_access.cpp](#)

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR

Structural node IR

Sparse
programming

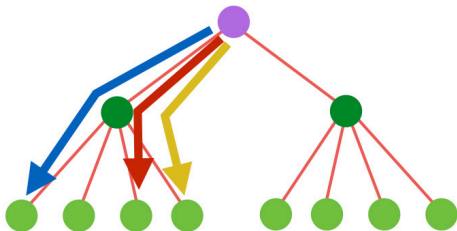
Differentiable
programming

Backends

Runtime system

Summary

Unoptimized Accesses



Optimized Accesses

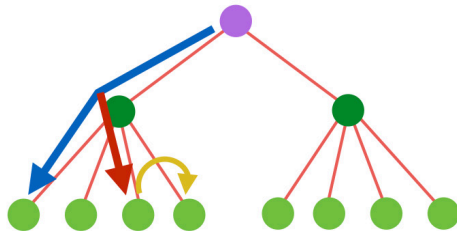


Figure: Access lowering breaks down end-to-end accesses into level-by-level micro-access operations. **Left:** An example unoptimized access; **Right:** After access lowering optimization passes can merge redundant data structure accesses.



Table of Contents

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

- 1 Introduction
- 2 Python frontend
- 3 Intermediate representation
 - Computation IR
 - Structural node IR
- 4 Sparse programming**
- 5 Differentiable programming
- 6 Backends
 - Runtime system
- 7 Summary

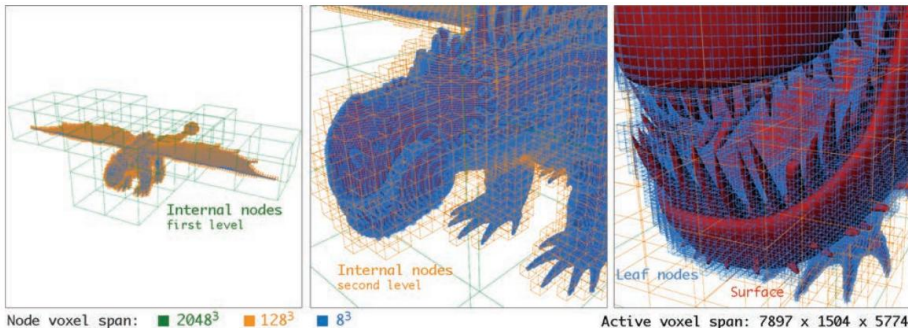


Fig. 4. High-resolution VDB created by converting polygonal model from *How To Train Your Dragon* to a narrow-band level set. The bounding resolution of the 228 million active voxels is $7897 \times 1504 \times 5774$ and the memory footprint of the VDB is 1GB, versus the $\frac{1}{4}$ TB for a corresponding dense volume. This VDB is configured with **LeafNodes** (blue) of size 8^3 and two levels of **InternalNodes** (green/orange) of size 16^3 . The index extents of the various nodes are shown as colored wireframes, and a polygonal mesh representation of the zero level set is shaded red. Images are courtesy of *DreamWorks Animation*.

²K. Museth (2013). “VDB: High-resolution sparse volumes with dynamic topology”. In: *ACM transactions on graphics (TOG)* 32.3, pp. 1–22.



OpenVDB³

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

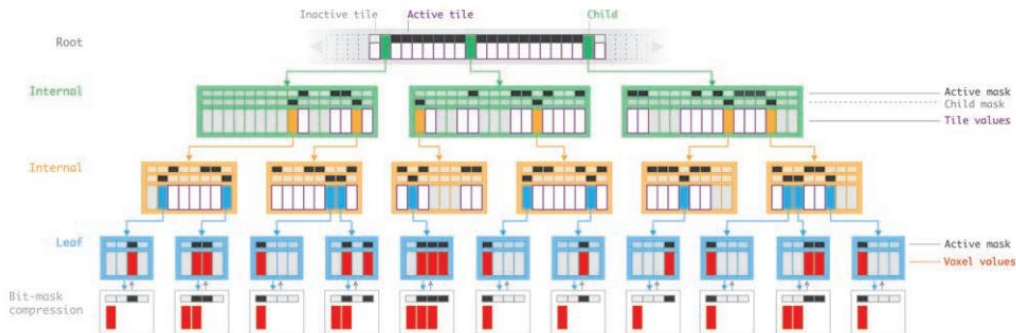
Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary



³K. Museth (2013). “VDB: High-resolution sparse volumes with dynamic topology”. In: *ACM transactions on graphics (TOG)* 32.3, pp. 1–22.



Spatially sparse programming

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

Ideally...

Save memory and computation on inactive spaces.

Reality...

Not easy to achieve desired performance and productivity

Taichi's solutions

- ① Allow programmers to flexibly define sparse data structures using SNodes
- ② Let programmers access sparse data structures as if they are dense
- ③ Automatically manage memory
- ④ Taichi's compiler automatically optimizes sparse data structure access



People's favorite sparse SNode

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

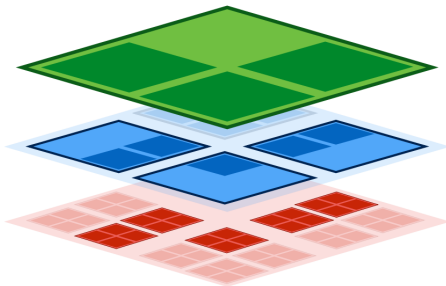
Differentiable
programming

Backends
Runtime system

Summary

The “pointer” SNode

- Essentially an array with each element being a pointer
- Different from the **dense** SNode, since pointers can be `nullptr`
- The most frequently used SNode to achieve sparsity



`ti.root`

`.pointer(ti.ij, 2)`

`.pointer(ti.ij, 2)`

`.dense(ti.ij, 2)`



“Pointer” and “Dense” are best buddies 😊

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

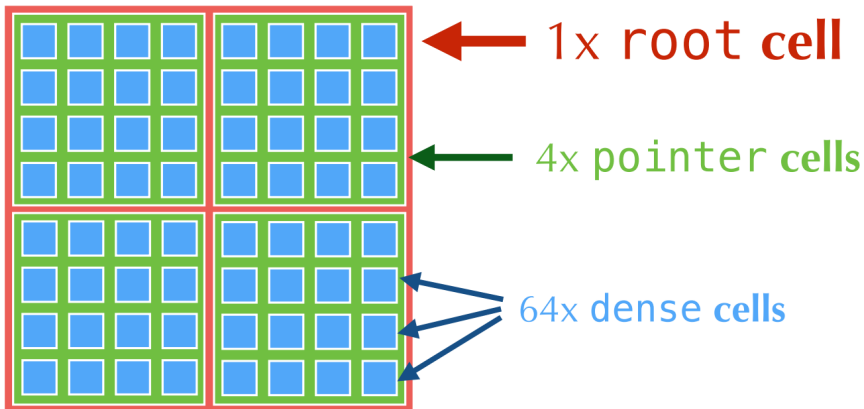
Differentiable
programming

Backends

Runtime system

Summary

```
block = ti.root.pointer(ti.ij, 2).dense(ti.ij, 4)
```





Sparse for-loops: no waste computation on inactive cells

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

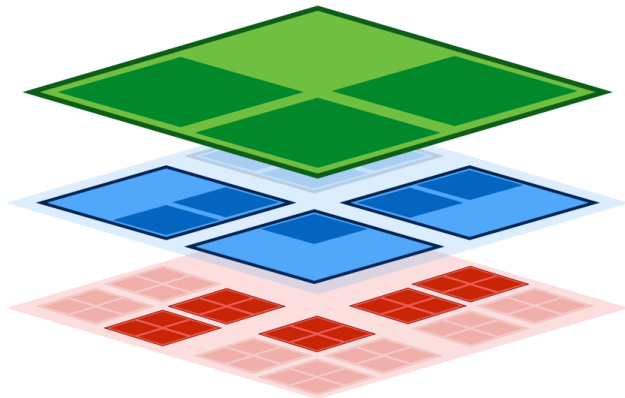
Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary



```
for i, j in a:  
    a[i, j] += 1
```




Table of Contents

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

- 1 Introduction
- 2 Python frontend
- 3 Intermediate representation
 - Computation IR
 - Structural node IR
- 4 Sparse programming
- 5 Differentiable programming**
- 6 Backends
 - Runtime system
- 7 Summary



Differentiable Programming

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

Forward programs evaluate $f(\mathbf{x})$; backward (gradient) programs evaluate $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$.

Taichi supports **reverse-mode automatic differentiation (AutoDiff)** that back-propagates gradients w.r.t. a scalar (loss) function $f(\mathbf{x})$.

Two ways to compute gradients:

- 1 Use Taichi's tape (`ti.Tape(loss)`) for both forward and gradient evaluation.
- 2 Explicitly use **gradient kernels** for gradient evaluation with more controls.



Gradient-based optimization

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

$$\min_{\mathbf{x}} \quad L(\mathbf{x}) = \frac{1}{2} \sum_{i=0}^{n-1} (\mathbf{x}_i - \mathbf{y}_i)^2.$$

- 1 Allocating fields with gradients:

```
x = ti.field(dtype=ti.f32, shape=n, needs_grad=True)
```

- 2 Defining loss function kernel(s):

```
@ti.kernel
def reduce():
    for i in range(n):
        L[None] += 0.5 * (x[i] - y[i])**2
```

- 3 Compute loss with `ti.Tape(loss=L): reduce()`

- 4 Gradient descent: `for i in x: x[i] -= x.grad[i] * 0.1`

Demo: `ti` example autodiff_minimization

Another demo: `ti` example autodiff_regression



Application 1: Forces from potential energy gradients

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

From the definition of potential energy:

$$\mathbf{f}_i = -\frac{\partial U(\mathbf{x})}{\partial \mathbf{x}_i}$$

Manually deriving gradients is hard. Let's use AutoDiff:

- 1 Allocate a 0D field to store the potential energy:
`potential = ti.field(ti.f32, shape=()).`
- 2 Define forward kernels that computes potential energy from `x[i]`.
- 3 In a `ti.Tape(loss=potential)`, call the forward kernels.
- 4 Force on each particle is `-x.grad[i]`.



Application 2: Differentiating a whole physical process

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

10 Demos: `DiffTaichi` $(\mathbf{x}_{t+1}, \mathbf{v}_{t+1}, \dots) = \mathbf{F}(\mathbf{x}_t, \mathbf{v}_t, \dots)$

Pattern:

```
with ti.Tape(loss=loss):  
    for i in range(steps - 1):  
        simulate(i)
```

Computational history

Always keep the whole computational history of time steps for end-to-end differentiation. I.e., instead of only allocating

`ti.Vector.field(3, dtype=ti.f32, shape=(num_particles))` that stores the latest particles, allocate for the whole simulation process

`ti.Vector.field(3, dtype=ti.f32, shape=(num_timesteps, num_particles))`. Do not overwrite! (Use **checkpointing** to reduce memory consumption.)



Flattening-based reverse-mode AD⁴

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

```
int a = 0;
if (b > 0) { a = b;}
else      { a = 2b;}
a = a + 1;
return a;
```

```
// flatten branching
int a = 0;
a = select(b > 0, b, 2b);
a = a + 1
return a;
```

```
// eliminate mutable var
ssa1 = select(b > 0, b, 2b);
ssa2 = ssa1 + 1
return ssa2;
```

⁴Y. Hu et al. (2020). “DiffTaichi: Differentiable Programming for Physical Simulation”. In: *ICLR*.



Flattening-based reverse-mode AD with “good” loops

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

forward

$$y_i = \sin x_i^2$$

```
for i in range(0, 16, step 1) do
    %1 = load x[i]
    %2 = mul %1, %1
    %3 = sin(%2)
    store y[i] = %3
end for
```

Backward

```
for i in range(0, 16, step 1) do
    // adjoint variables
    %1adj = alloca 0.0
    %2adj = alloca 0.0
    %3adj = alloca 0.0
    // original forward computation
    %1 = load x[i]
    %2 = mul %1, %1
    %3 = sin(%2)
    // reverse accumulation
    %4 = load y_adj[i]
    %3adj += %4
    %5 = cos(%2)
    %2adj += %3adj * %5
    %1adj += 2 * %1 * %2adj
    atomic add x_adj[i], %1adj
end for
```

Note: Taichi uses global tensors as checkpoints,
so we need some recomputation.
(Megakernel has a cost in AutoDiff but still helps.)



Stack-based reverse-mode AD

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

Slightly harder cases where we need stacks to trace local computational history:

► [transforms/auto_diff.cpp](#)

```
@ti.kernel
def fib():
    for i in range(N):
        p = a[i]
        q = b[i]
        for j in range(c[i]):
            new_p = q
            new_q = p + q
            p, q = new_p, new_q
        f[i] = q
```

```
@ti.kernel
def power():
    for i in range(N):
        ret = 1.0
        for j in range(b[i]):
            ret = ret * a[i]
        p[i] = ret
```




Stack-based reverse-mode AD (Control flow)

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

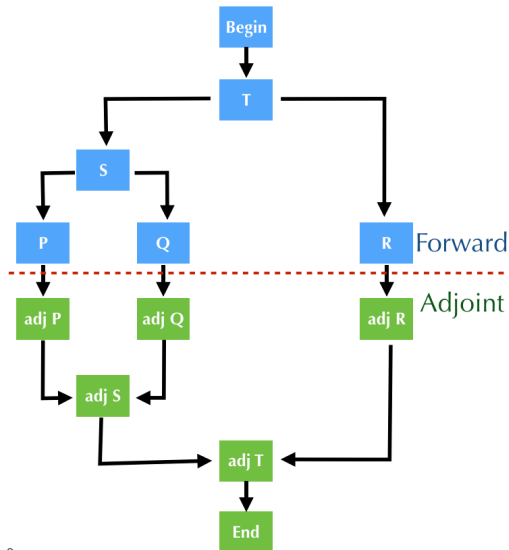
Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

```
T
if (x) {
  S
  if (y) {
    P
  } else {
    Q
  }
} else {
  R
}
```





Two-scale AutoDiff ▶ lang/tape.py

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

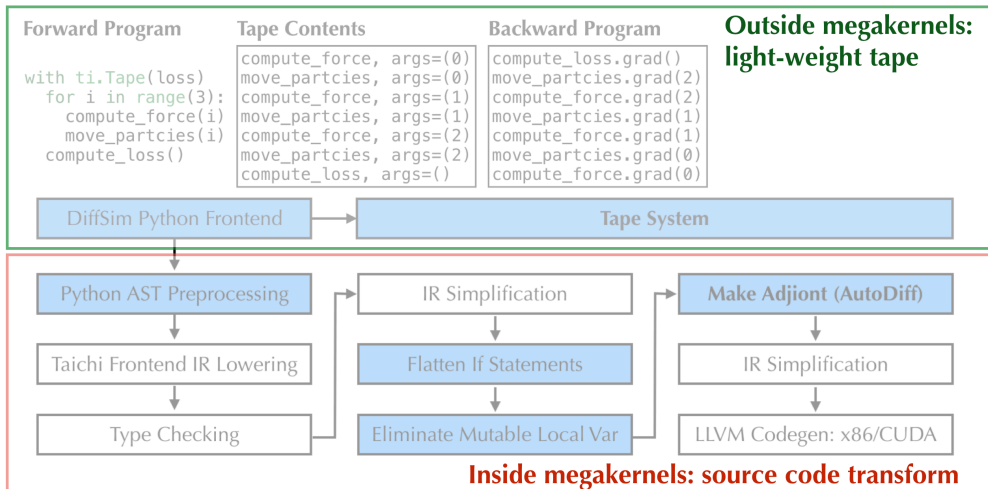




Table of Contents

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

- 1 Introduction
- 2 Python frontend
- 3 Intermediate representation
 - Computation IR
 - Structural node IR
- 4 Sparse programming
- 5 Differentiable programming
- 6 Backends**
 - Runtime system
- 7 Summary



7 existing backends of Taichi

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

Taichi backends ▶ [backends](#)

LLVM ▶ [codegen/codegen_llvm.cpp](#)

CPU: x64, AArch64 (aka ARM64)

▶ [backends/cpu](#)

CUDA (NVPTX) ▶ [backends/cuda](#)

C99 ▶ [backends/cc](#)

JavaScript/WebAssembly: [Demos](#)

OpenGL ▶ [backends/opengl](#)

Metal ▶ [backends/metal](#)



Describing backend capabilities

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

Not all backends are fully featured

Due to backend API limitations it may not be possible to implement certain features on all backends. For example, Apple Metal doesn't support 64-bit data types, so `ti.f64` is not available on Taichi's Metal backend.

Backend extension system

Backend extensions are used to describe the capability of each backend (beyond core functionalities). For example,

- 1 `Extension::sparse:` sparse computation
- 2 `Extension::data64:` 64-bit data types
- 3 `Extension::adstack:` stack-based AutoDiff
- 4 All backend extensions: [▶inc/extensions.inc.h](#)

Which backend supports what? [▶program/extension.cpp](#)



The LLVM backends

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

The CPU and CUDA code generators are based on LLVM (v10.0.0) ORCv2 JIT (CPU) and NVPTX backend (CUDA).

- Extremely portable: the LLVM CPU backend has no dependency
- Used as the fallback solution for all other backends
- Fully featured (i.e., support all extensions)
- Serve as references for other backends
- Has a (ideally) zero-cost runtime system (more on this later)

▶ [runtime/llvm/runtime.cpp](#)



Source-to-source backends

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

Most backends do not take LLVM IR as input. In those cases we directly emit source code (i.e., source-to-source compilation).

- ① C99 ▶ [backends/cc](#)
 - + Emscripten = [JavaScript/WebAssembly](#)
- ② OpenGL ▶ [backends/opengl](#)
- ③ Metal ▶ [backends/metal](#)

Legacy CPU and CUDA backends

Taichi used to emit C++ and CUDA source code for CPU and CUDA backends as well. However, this solution is abandoned due to portability issues and long compilation time.



LLVM runtime system [▶runtime/llvm](#)

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

Not really a runtime system...

Design tradeoffs

- One implementation should ideally serve for all LLVM-related backends
- The runtime code need to be compiled into LLVM IR and then linked against generated code
- We want to use something like C++ templates for each SNode
- ... but we cannot do “link-time template specialization” on LLVM IR \implies Inline as much as possible to “simulate” templates.

Runtime system written in C++ [▶runtime/llvm/runtime.cpp](#) will be compiled to LLVM IR using `clang` by developers and ship together with the Taichi package.



Memory allocator: design constraints

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

- ① Portability matters
 - Need to support devices with/without (unified) virtual memory
- ② Need high-performance, in parallel
 - Make it as lock-free as possible
 - Trade virtual address space and external fragmentation for performance
- ③ Need to support many backends
 - Make it as simple as possible
 - Use (mostly) the same code path for CPU/GPUs



Memory allocator: the “big” picture

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

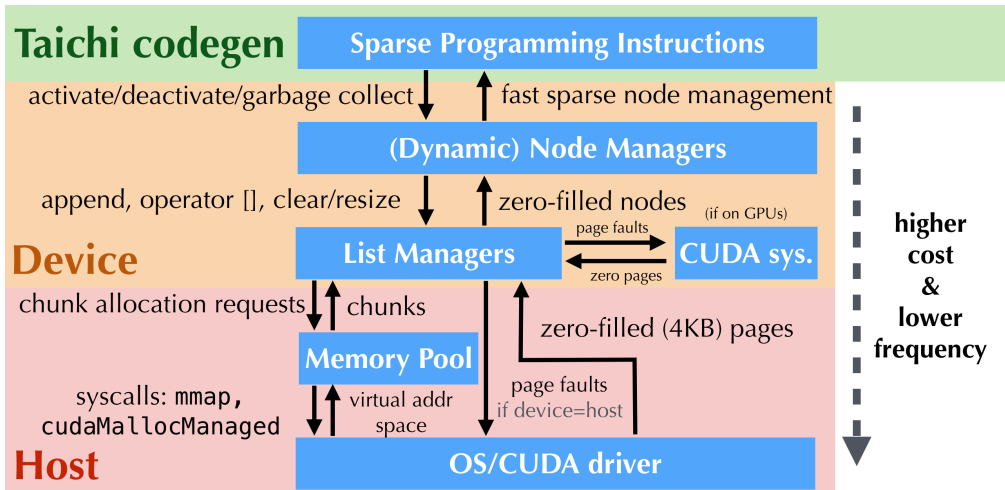




Table of Contents

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends
Runtime system

Summary

- 1 Introduction
- 2 Python frontend
- 3 Intermediate representation
 - Computation IR
 - Structural node IR
- 4 Sparse programming
- 5 Differentiable programming
- 6 Backends
 - Runtime system
- 7 Summary



Summary

Life of a Taichi
Kernel

Yuanming Hu

Introduction

Python frontend

Intermediate
representation

Computation IR
Structural node IR

Sparse
programming

Differentiable
programming

Backends

Runtime system

Summary

- Language design and compiler engineering is a lot of fun!
- Keep it simple and practical.
- **Portability** is not an easy task: took me 4 months to get a prototype (and a SIGGRAPH Asia 2019 paper), but more than 1 year to make it portable.
- ~ 400,000 downloads according to PePy, since Jan 2020.

Note

Many features of Taichi are developed by **the Taichi community**.
Clearly, I am not the only developer 😊

More details...

- [GitHub](#)
- [Docs](#)
- [TaichiCon](#)
- [SIGGRAPH 2020 Taichi course](#)