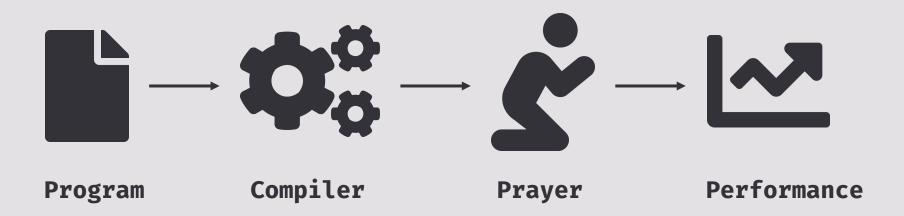
# ELEV//TE

### A Language for Expressing Optimization Strategies

Bastian Hagedorn (University of Münster) Martin Lücke (University of Münster) Johannes Lenfers (University of Münster) Michel Steuwer (University of Glasgow)

How do we optimize programs today?



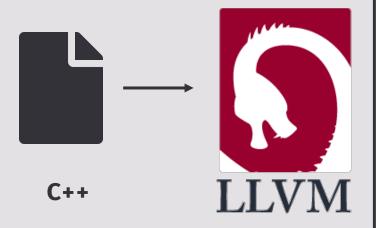
How do we optimize programs today?

### **General Purpose Compilers**



#### How do we optimize programs today?

### **General Purpose Compilers**

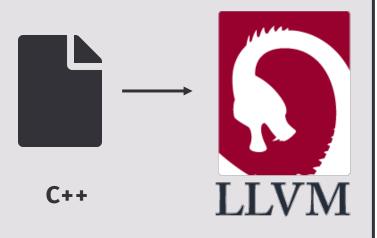


```
Code Generation Options
-00, -01, -02, -03, -0fast, -0s, -0z, -0g, -0, -04
    Specify which optimization level to use:
          -00 Means "no optimization": this level compiles the fastest and generates the most debuggable code.
          -01 Somewhere between -00 and -02.
          -02 Moderate level of optimization which enables most optimizations.
          -03 Like -02, except that it enables optimizations that take longer to perform or that may generate larger code (in an
          attempt to make the program run faster).
          -Ofast Enables all the optimizations from -03 along with other aggressive optimizations that may violate strict
          compliance with language standards.
          -0s Like -02 with extra optimizations to reduce code size.
          -0z Like -0s (and thus -02), but reduces code size further.
          -og Like -o1. In future versions, this option might disable different optimizations in order to improve debuggability.
          -0 Equivalent to -02.
          -04 and higher
                Currently equivalent to -03
```

### Code Generation Options

#### How do we optimize programs today?

### **General Purpose Compilers**

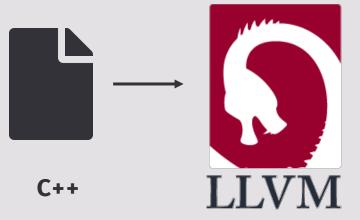


```
Code Generation Options
-00, -01, -02, -03, -0fast, -0s, -0z, -0g, -0, -04
    Specify which optimization level to use:
          -00 Means "no optimization": this level compiles the fastest and generates the most debuggable code.
          -01 Somewhere between -00 and -02.
          -02 Moderate level of optimization which enables most optimizations.
          -03 Like -02, except that it enables optimizations that take longer to perform or that may generate larger code (in an
          attempt to make the program run faster).
          -Ofast Enables at the optimizations from -03 along with other aggressive optimizations that may violate strict
          compliance with language standards.
          -0s Like -02 with extra optimizations to reduce code size.
          -0z Like -0 (and thus -02), but reduces code size further.
          -og Like of. In future versions, this option might disable different optimizations in order to improve debuggability.
          -0 Equivalent to -02.
             and higher
                Currently equivalent to -03
```

"... in an attempt to make the program run faster"

#### How do we optimize programs today?

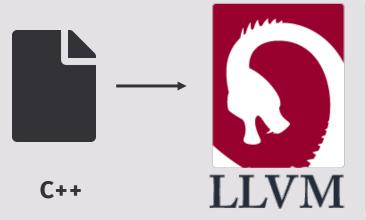
### **General Purpose Compilers**



-targetlibinfo -tti -tbaa -scoped-noalias -assumption-cache-tracker -profile-summary-info -forceattrs -inferattrs -callsite-splitting psccp -called-value-propagation -globalopt -domtree -mem2reg -deadargelim -domtree -basicaa -aa -loops -lazy-branch-prob -lazy-block-fre q -opt-remark-emitter -instcombine -simplifycfg -basiccg -globals-aa -prune-eh -inline -functionattrs -argpromotion -domtree -sroa -basi caa -aa -memoryssa -early-cse-memssa -speculative-execution -basicaa -aa -lazy-value-info -jump-threading -correlated-propagation -simpl ifycfg -domtree -aggressive-instcombine -basicaa -aa -loops -lazy-branch-prob -lazy-block-freq -opt-remark-emitter -instcombine -libcal s-shrinkwrap -loops -branch-prob -block-freq -lazy-branch-prob -lazy-block-freq -opt-remark-emitter -pgo-memop-opt -basicaa -aa -loops lazy-branch-prob -lazy-block-freq -opt-remark-emitter -tailcallelim -simplifycfg -reassociate -domtree -loops -loop-simplify -lcssa-ver fication -lcssa -basicaa -aa -scalar-evolution -loop-rotate -licm -loop-unswitch -simplifycfg -domtree -basicaa -aa -loops -lazy-branch prob -lazv-block-freg -opt-remark-emitter -instcombine -loop-simplify -lcssa-verification -lcssa -scalar-evolution -indvars -loop-idiom -loop-deletion -loop-unroll -mldst-motion -phi-values -basicaa -aa -memdep -lazy-branch-prob -lazy-block-freq -opt-remark-emitter -gvn phi-values -basicaa -aa -memdep -memcpyopt -sccp -demanded-bits -bdce -basicaa -aa -loops -lazy-branch-prob -lazy-block-freq -opt-remark -emitter -instcombine -lazy-value-info -jump-threading -correlated-propagation -basicaa -aa -phi-values -memdep -dse -loops -loop-simpli fy -lcssa-verification -lcssa -basicaa -aa -scalar-evolution -licm -postdomtree -adce -simplifycfg -domtree -basicaa -aa -loops -lazy-br anch-prob -lazv-block-freq -opt-remark-emitter -instcombine -barrier -elim-avail-extern -basiccg -rpo-functionattrs -globalopt -globaldc e -basiccg -globals-aa -float2int -domtree -loops -loop-simplify -lcssa-verification -lcssa -basicaa -aa -scalar-evolution -loop-rotate -loop-accesses -lazy-branch-prob -lazy-block-freq -opt-remark-emitter -loop-distribute -branch-prob -block-freq -scalar-evolution -basi aa -aa -loop-accesses -demanded-bits -lazy-branch-prob -lazy-block-freq -opt-remark-emitter -loop-vectorize -loop-simplify -scalar-evolu tion -aa -loop-accesses -loop-load-elim -basicaa -aa -lazy-branch-prob -lazy-block-freq -opt-remark-emitter -instcombine -simplifycfg omtree -loops -scalar-evolution -basicaa -aa -demanded-bits -lazy-branch-prob -lazy-block-freq -opt-remark-emitter -slp-vectorizer -optremark-emitter -instcombine -loop-simplify -lcssa-verification -lcssa -scalar-evolution -loop-unroll -lazy-branch-prob -lazy-block-freq opt-remark-emitter -instcombine -loop-simplify -lcssa-verification -lcssa -scalar-evolution -licm -alignment-from-assumptions -strip-de ad-prototypes -globaldce -constmerge -domtree -loops -branch-prob -block-freq -loop-simplify -lcssa-verification -lcssa -basicaa -aa -so alar-evolution -branch-prob -block-freq -loop-sink -lazy-branch-prob -lazy-block-freq -opt-remark-emitter -instsimplify -div-rem-pairs simplifycfg -verify

#### How do we optimize programs today?

### **General Purpose Compilers**



targetlibinfo -tti -tbaa -scoped-noalias -assumption-cache-tracker -profile-summa sccp -called-value-propagation -globalopt -domtree -mem2reg -deadargelim -domtree -opt-remark-emitter -instcombine -simplifycfg -basiccg -globals-aa -prune-eh -inl aa -aa -memoryssa -early-cse-memssa -speculative-execution -basicaa -aa -lazy-valu ifycfg -domtree -aggressive-instcombine -basicaa -aa -loops -lazy-branch-prob -lazy -shrinkwrap -loops -branch-prob -block-freq -lazy-branch-prob -lazy-block-freq -op azy-branch-prob -lazy-block-freq -opt-remark-emitter -tailcallelim -simplifycfg ication -lcssa -basicaa -aa -scalar-evolution -loop-rotate -licm -loop-unswitch prob -lazy-block-freq -opt-remark-emitter -instcombine -loop-simplify -lcssa-verif loop-deletion -loop-unroll -mldst-motion -phi-values -basicaa -aa -memdep -lazy-br hi-values -basicaa -aa -memdep -memcpyopt -sccp -demanded-bits -bdce -basicaa -aa emitter -instcombine -lazy-value-info -jump-threading -correlated-propagation -bas v -lcssa-verification -lcssa -basicaa -aa -scalar-evolution -licm -postdomtree -ac nch-prob -lazy-block-freq -opt-remark-emitter -instcombine -barrier -elim-avail-ex -basiccg -globals-aa -float2int -domtree -loops -loop-simplify -lcssa-verificatio loop-accesses -lazy-branch-prob -lazy-block-freq -opt-remark-emitter -loop-distri a -aa -loop-accesses -demanded-bits -lazy-branch-prob -lazy-block-freq -opt-remark ion -aa -loop-accesses -loop-load-elim -basicaa -aa -lazy-branch-prob -lazy-blockmtree -loops -scalar-evolution -basicaa -aa -demanded-bits -lazy-branch-prob -lazy remark-emitter -instcombine -loop-simplify -lcssa-verification -lcssa -scalar-evolu opt-remark-emitter -instcombine -loop-simplify -lcssa-verification -lcssa -scalarad-prototypes -globaldce -constmerge -domtree -loops -branch-prob -block-freq -loop lar-evolution -branch-prob -block-freq -loop-sink -lazy-branch-prob -lazy-block-f

lazv-block-f

a -aa -loops

fv -lcssa-ver

-lazy-branch

s -loop-idion

-emitter -gvn

eq -opt-remai

s -loop-simpl

loops -lazy-b

lopt -globalo

-loop-rotate

olution -basi

-scalar-evol

-simplifycfg

ctorizer -op

zy-block-freq

ions -strip-d

basicaa -aa -s

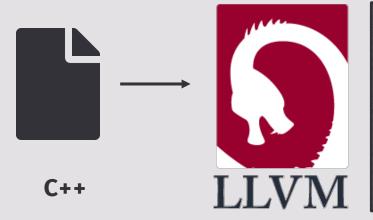
div-rem-pairs

Compiler

Passes

#### How do we optimize programs today?

### **General Purpose Compilers**



targetlibinfo -tti -tbaa -scoped-noalias -assumption-cache-tracker -profile-summa sccp -called-value-propagation -globalopt -domtree -mem2reg -deadargelim -domtree -opt-remark-emitter -instcombine -simplifycfg -basiccg -globals-aa -prune-eh -inl aa -aa -memoryssa -early-cse-memssa -speculative-execution -basicaa -aa -lazy-valu lfycfg -domtree -aggressive-<mark>instcombine</mark> -basicaa -aa -loops -lazy-branch-prob -lazy -shrinkwrap -loops -branch-prob -block-freq -lazy-branch-prob -lazy-block-freq -o azy-branch-prob -lazy-block-freq -opt-remark-emitter -tailcallelim -simplifycfg fication -lcssa -basicaa -aa -scalar-evolution -loop-rotate -licm -loop-unswitch rob -lazy-block-freq -opt-remark-emitter -<mark>instcombine</mark> -loop-simplify -lcssa-verif: -loop-deletion -loop-unroll -mldst-motion -phi-values -basicaa -aa -memdep -lazy-br hi-values -basicaa -aa -memdep -memcpyopt -sccp -demanded-bits -bdce -basicaa -aa emitter -<mark>instcombine</mark> -lazy-value-info -jump-threading -correlated-propagation -bas v -lcssa-verification -lcssa -basicaa -aa -scalar-evolution -licm -postdomtree -ad nch-prob -lazy-block-freq -opt-remark-emitter -<mark>instcombine</mark> -barrier -elim-avail-ex -basiccg -globals-aa -float2int -domtree -loops -loop-simplify -lcssa-verification loop-accesses -lazy-branch-prob -lazy-block-freq -opt-remark-emitter -loop-distrib a -aa -loop-accesses -demanded-bits -lazy-branch-prob -lazy-block-freq -opt-remark ion -aa -loop-accesses -loop-load-elim -basicaa -aa -lazy-branch-prob -lazy-block mtree -loops -scalar-evolution -basicaa -aa -demanded-bits -lazy-branch-prob -lazy remark-emitter -<mark>instcombine</mark> -loop-simplify -lcssa-verification -lcssa -scalar-evolu opt-remark-emitter -<mark>instcombine</mark> -loop-simplify -lcssa-verification -lcssa -scalar d-prototypes -globaldce -constmerge -domtree -loops -branch-prob -block-freq -loop lar-evolution -branch-prob -block-freq -loop-sink -lazy-branch-prob -lazy-block-f

Institution could be a served on the country of the country o

lazv-block-fi ee -sroa -bas agation -sim nbine -libcal a -aa -loops fy -lcssa-ver -lazy-branch s -loop-idion -emitter -gvn req -opt-remai s -loop-simpl -loops -lazy-b alopt -globalo -loop-rotate olution -bas -scalar-evol -simplifycfg ectorizer -opt azv-block-fred tions -strip-o basicaa -aa -s div-rem-pairs

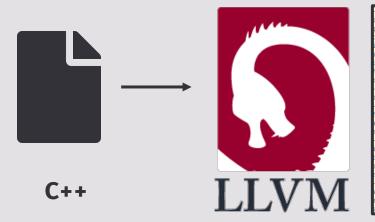
-splitting

-03

Compiler Passes

#### How do we optimize programs today?

### **General Purpose Compilers**



targetlibinfo -tti -tbaa -scoped-noalias -assumption-cache-tracker -profile-summa sccp -called-value-propagation -globalopt -domtree -mem2reg -deadargelim -domtree -opt-remark-emitter -instcombine -simplifycfg -basiccg -globals-aa -prune-eh -inl aa -aa -memoryssa -early-cse-memssa -speculative-execution -basicaa -aa -lazy-valu lfycfg -domtree -aggressive-<mark>instcombine</mark> -basicaa -aa -loops -lazy-branch-prob -lazy -shrinkwrap -loops -branch-prob -block-freq -lazy-branch-prob -lazy-block-freq -op azy-branch-prob -lazy-block-freq -opt-remark-emitter -tailcallelim -simplifycfg fication -lcssa -basicaa -aa -scalar-evolution -loop-rotate -licm -loop-unswitch rob -lazy-block-freq -opt-remark-emitter -<mark>instcombine</mark> -loop-simplify -lcssa-verif: -loop-deletion -loop-unroll -mldst-motion -phi-values -basicaa -aa -memdep -lazy-br hi-values -basicaa -aa -memdep -memcpyopt -sccp -demanded-bits -bdce -basicaa -aa emitter -<mark>instcombine</mark> -lazy-value-info -jump-threading -correlated-propagation -bas v -lcssa-verification -lcssa -basicaa -aa -scalar-evolution -licm -postdomtree -ac nch-prob -lazy-block-freq -opt-remark-emitter -instcombine -barrier -elim-avail-ex -basiccg -globals-aa -float2int -domtree -loops -loop-simplify -lcssa-verification loop-accesses -lazy-branch-prob -lazy-block-freq -opt-remark-emitter -loop-distrib a -aa -loop-accesses -demanded-bits -lazy-branch-prob -lazy-block-freq -opt-remark ion -aa -loop-accesses -loop-load-elim -basicaa -aa -lazy-branch-prob -lazy-block mtree -loops -scalar-evolution -basicaa -aa -demanded-bits -lazy-branch-prob -lazy remark-emitter -<mark>instcombine</mark> -loop-simplify -lcssa-verification -lcssa -scalar-evolu opt-remark-emitter -instcombine -loop-simplify -lcssa-verification -lcssa -scalard-prototypes -globaldce -constmerge -domtree -loops -branch-prob -block-freq -loop alar-evolution -branch-prob -block-freq -loop-sink -lazy-branch-prob -lazy-block-f

Controlling optimisations is hard, because: We have no clue what's going on inside the compiler!

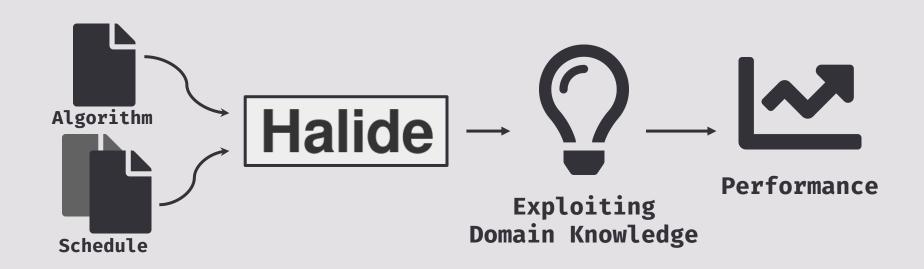
\*\* The control of the control o

.azv-block-fi ee -sroa -bas fy -lcssa-ver -lazy-branch s -loop-idio -emitter -gvn s -loop-simpl -loops -lazy-b alopt -globalo olution -bas: -scalar-evol -simplifycfg ectorizer -opt azv-block-fred tions -strip-o basicaa -aa -s div-rem-pairs

Compiler Passes

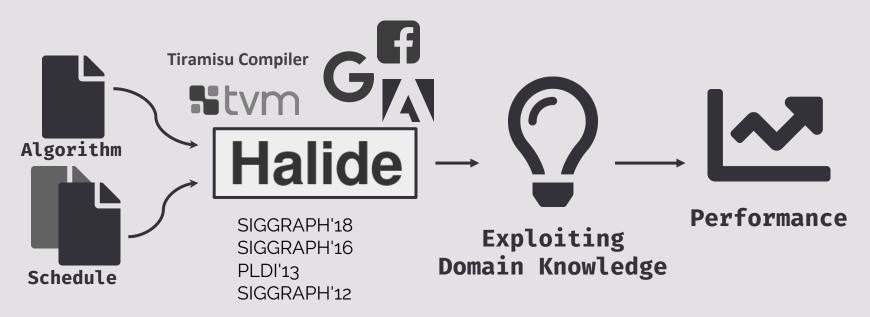
How do we optimize programs today?

### **Domain-Specific Compilers**



How do we optimize programs today?

### **Domain-Specific Compilers**



#### How do we optimize programs today?

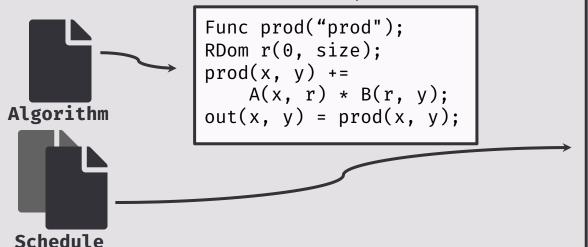
### **Domain-Specific Compilers**

Matrix Multiplication

How do we optimize programs today?

### **Domain-Specific Compilers**

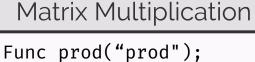
Matrix Multiplication

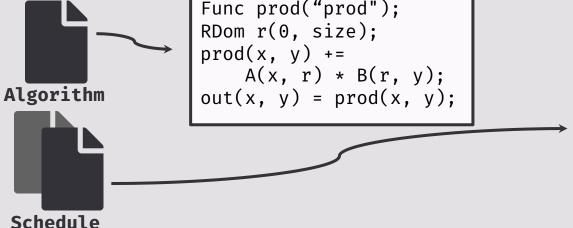


```
const int warp size = 32;
const int vec size = 2;
const int x tile = 3;
const int v tile = 4;
const int y_unroll = 8;
const int r unroll = 1;
Var xi, yi, xio, xii, yii, xo, yo, x_pair, xiio, ty;
RVar rxo, rxi;
out.bound(x, 0, size)
   .bound(v, 0, size)
   .tile(x, y, xi, yi, x_tile * vec_size * warp_size, y_tile *
   .split(vi, tv, vi, v unroll)
   .vectorize(xi, vec size)
   .split(xi, xio, xii, warp size)
   .reorder(xio, vi, xii, tv, x, v)
   .unroll(xio)
   .unroll(vi)
   .gpu blocks(x, y)
   .gpu threads(ty)
   .gpu_lanes(xii);
prod.store in(MemoryType!::Register)
                        .compute at(out, x)
                        .split(x, xo, xi, warp size * vec size,
TailStrategy!::RoundUp)
                        .split(y, ty, y, y_unroll)
                        .gpu threads(ty)
                        .unroll(xi, vec size)
                        .gpu lanes(xi)
                        .unroll(xo)
                        .unroll(v)
                        .update()
                        .split(x, xo, xi, warp size * vec size,
TailStrategy!::RoundUp)
                        .split(y, ty, y, y_unroll)
                        .gpu threads(ty)
                        .unroll(xi, vec size)
                        .gpu lanes(xi)
                        .split(r.x, rxo, rxi, warp_size)
                        .unroll(rxi, r unroll)
                        .reorder(xi, xo, v, rxi, tv, rxo)
                        .unroll(xo)
                        .unroll(v);
Var Bx = B.in().args()[0], By = B.in().args()[1];
Var Ax = A.in().args()[0], Av = A.in().args()[1];
B.in()
                        .compute at(prod, ty)
                        .split(Bx, xo, xi, warp_size)
                        .gpu lanes(xi)
                        .unroll(xo).unroll(By);
A.in()
                        .compute at(prod, rxo)
                        .vectorize(Ax, vec size)
                        .split(Ax, xo, xi, warp size)
                        .gpu lanes(xi)
                        .unroll(xo).split(Ay, yo, yi, y tile)
                        .gpu threads(vi).unroll(vo);
A.in().in().compute at(prod, rxi)
                        .vectorize(Ax, vec size)
                        .split(Ax, xo, xi, warp size)
                        .gpu lanes(xi)
                        .unroll(xo).unroll(Ay);
set alignment and bounds(A, size);
set alignment and bounds(B, size);
set alignment and bounds(out, size);
```

How do we optimize programs today?

### **Domain-Specific Compilers**



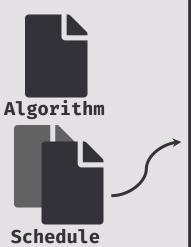


Schedules are much harder to write than algorithms

```
const int warp size = 32;
const int vec size = 2;
const int x tile = 3;
const int v tile = 4;
const int y unroll = 8;
const int r unroll = 1;
Var xi, yi, xio, xii, yii, xo, yo, x_pair, xiio, ty;
RVar rxo, rxi;
out.bound(x, 0, size)
   .bound(v, 0, size)
   .tile(x, y, xi, yi, x tile * vec size * warp size, y tile *
   .split(vi, tv, vi, v unroll)
   .vectorize(xi, vec size)
   .split(xi, xio, xii, warp size)
   .reorder(xio, vi, xii, tv, x, v)
   .unroll(xio)
   .unroll(vi)
   .gpu blocks(x, y)
   .gpu threads(ty)
   .gpu_lanes(xii);
prod.store in(MemoryType!::Register)
                        .compute at(out, x)
                        .split(x, xo, xi, warp size * vec size,
TailStrategy!::RoundUp)
                        .split(y, ty, y, y_unroll)
                        .gpu threads(ty)
                        .unroll(xi, vec size)
                        .gpu lanes(xi)
                        .unroll(xo)
                        .unroll(v)
                        .update()
                        .split(x, xo, xi, warp size * vec size,
TailStrategy!::RoundUp)
                        .split(y, ty, y, y_unroll)
                        .gpu threads(ty)
                        .unroll(xi, vec size)
                        .gpu lanes(xi)
                        .split(r.x, rxo, rxi, warp_size)
                        .unroll(rxi, r unroll)
                        .reorder(xi, xo, v, rxi, tv, rxo)
                        .unroll(xo)
                        .unroll(v);
Var Bx = B.in().args()[0], By = B.in().args()[1];
Var Ax = A.in().args()[0], Av = A.in().args()[1];
B.in()
                        .compute at(prod, ty)
                        .split(Bx, xo, xi, warp_size)
                        .gpu lanes(xi)
                        .unroll(xo).unroll(By);
A.in()
                        .compute at(prod, rxo)
                        .vectorize(Ax, vec size)
                        .split(Ax, xo, xi, warp size)
                        .gpu lanes(xi)
                        .unroll(xo).split(Ay, yo, yi, y tile)
                        .gpu threads(vi).unroll(vo);
A.in().in().compute at(prod, rxi)
                        .vectorize(Ax, vec size)
                        .split(Ax, xo, xi, warp size)
                        .gpu lanes(xi)
                        .unroll(xo).unroll(Ay);
set alignment and bounds(A, size);
set alignment and bounds(B, size);
set alignment and bounds(out, size);
```

#### How do we optimize programs today?

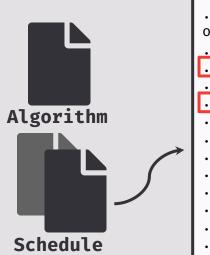
### **Domain-Specific Compilers**



```
out.bound(x, 0, size)
.bound(y, 0, size)
.tile(x, y, xi, yi, x_tile * vec_size * warp_size, y_tile * y_unroll)
.split(yi, ty, yi, y_unroll)
.vectorize(xi, vec size)
.split(xi, xio, xii, warp_size)
.reorder(xio, yi, xii, ty, x, y)
.unroll(xio)
.unroll(yi)
.gpu blocks(x, y)
.gpu threads(ty)
.gpu lanes(xii);
```

How do we optimize programs today?

Domain-Specific Compilers fixed set of optimizations  $\Rightarrow$  lack of extensibility

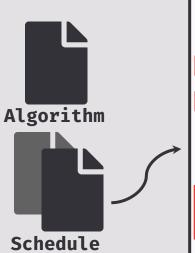


```
out.bound(x, 2, size)
.bound(y 0, size)
.tile(x, y, xi, yi, x_tile * vec_size * warp_size, y_tile * y_unroll)
.split(yi, ty, yi, y_unroll)
.vectorize(ki, vec size)
.split(xi, xio, xii, warp_size)
.reorder(xio, yi, xii, ty, x, y)
.unroll(xio)
.unroll(yi)
.gpu blocks(x, y)
.gpu threads(ty)
.gpu lanes(xii);
```

How do we optimize programs today?

**Domain-Specific Compilers** 

fixed set of optimizations  $\Rightarrow$  lack of extensibility

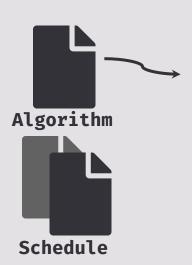


```
out.bound(x, 9, size)
.bound(v 0, size)
.tile(x, y, xi, yi, x_tile * vec_size * warp_size, y_tile * y unroll)
.split(yi, ty, yi, y_unroll)
.vectorize(ki, vec size)
.split(xi, xio, xii, warp_size)
.reorder(xio, yi, xii, ty, x, y)
.unroll(xio)
.unroll(yi)
.gpu_blocks(x, y) what happens if the order of these is swapped?
.gpu threads(ty)
                 ⇒ unclear semantics
.gpu_lanes(xii);
                 ⇒ unclear how to automatically generate schedules
```

#### How do we optimize programs today?

Not always a clear separation of algorithm and schedule

**Domain-Specific Compilers** 

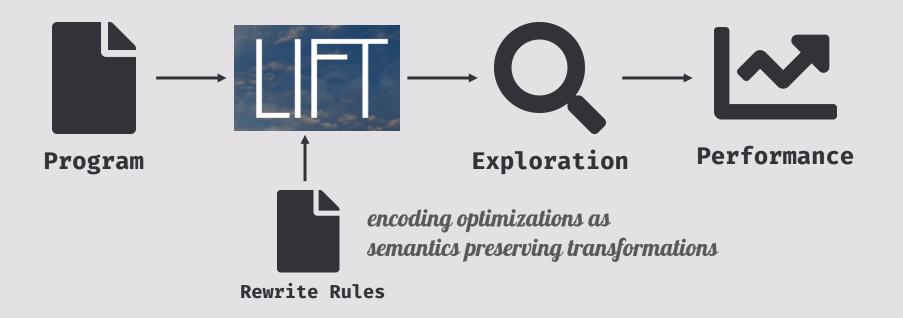


```
Func prod("prod");
prod(j, i) = A_{(j, i)} * x_{(j)};
RDom k(0, sum_size_vecs, "k");
Func accum_vecs("accum_vecs"):
accum_vecs(j, i) += prod(k * vec_size + j, i);
Func accum vecs transpose("accum vecs transpose");
accum vecs transpose(i, j) = accum vecs(j, i);
RDom lanes(0, vec_size);
Func sum_lanes("sum_lanes");
sum lanes(i) += accum vecs transpose(i, lanes);
RDom tail(sum_size_vecs * vec_size, sum_size - sum_size_vecs * vec_size);
Func sum_tail("sum_tail");
sum_tail(i) = sum_lanes(i);
sum tail(i) += prod(tail, i);
Func Ax("Ax");
Ax(i) = sum tail(i);
result(i) = b * v (i) + a * Ax(i);
```

Matrix Vector Multiplication

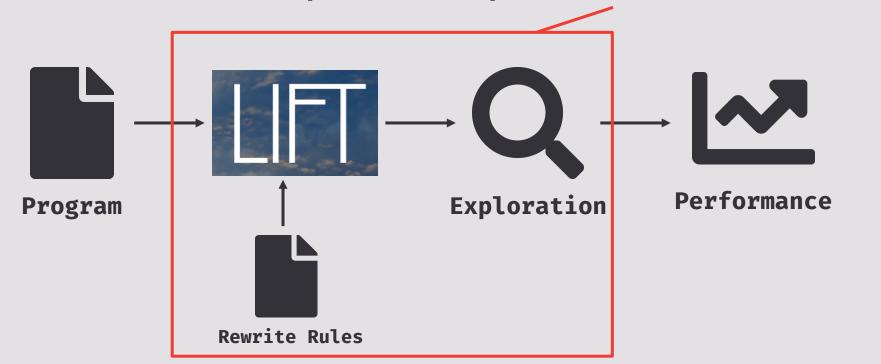
How do we optimize programs today?

### **Functional Domain-Specific Compilers**



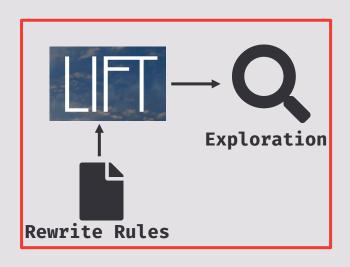
How do we optimize programs today?

Functional Domain-Specific Compilers How does this actually work?



How do we optimize programs today?

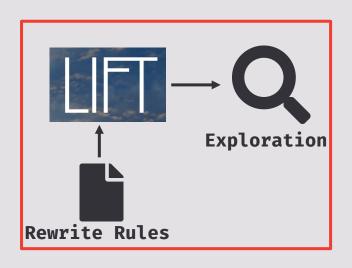
### **Functional Domain-Specific Compilers**



CGO'18 CGO'17 CASES'16 GPGPU'16 ICFP'15

How do we optimize programs today?

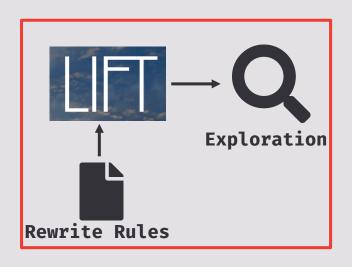
### **Functional Domain-Specific Compilers**



CGO'18 No explanation CGO'17 CASES'16 GPGPU'16 ICFP'15

How do we optimize programs today?

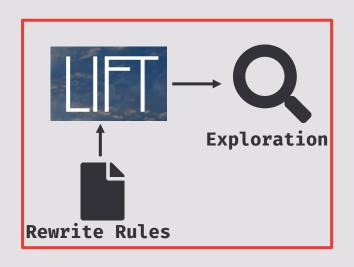
### **Functional Domain-Specific Compilers**



CGO'18
CGO'17 No explanation
CASES'16
GPGPU'16
ICFP'15

How do we optimize programs today?

### **Functional Domain-Specific Compilers**



CGO'18 CGO'17 CASES'16 CASES'16 CASES'16 CASES'16 ...the resulting space is extremely large, even potentially unbounded, which opens up a *new research challenge*.

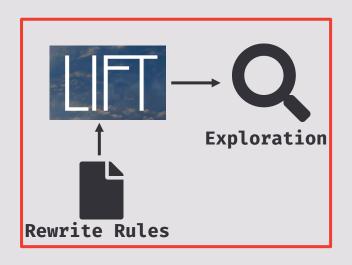
•••

We present here a first, *simple and heuristic-based* pruning strategy to tackle the space complexity problem. *Future research* will investigate more advanced techniques to fully automate the pruning process.

. . .

How do we optimize programs today?

### **Functional Domain-Specific Compilers**



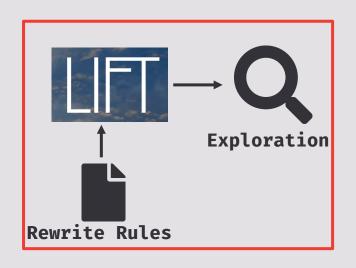
CGO'18 CGO'17 CASES'16 GPGPU'16 ICFP'15 ...guide the automatic rewrite process by grouping rewrite rules together into macro rules

...

These macro rules are more flexible than the simple rules. They *try to apply different sequences* of rewrites to achieve their optimization goal.

How do we optimize programs today?

### **Functional Domain-Specific Compilers**



CGO'18 CGO'17 CASES'16 GPGPU'16 ICFP'15 ...guide the automatic rewrite process by grouping rewrite rules together into macro rules

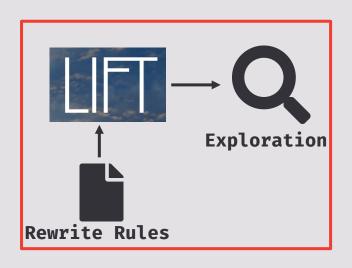
...

These macro rules are more flexible than the simple rules. They *try to apply different sequences* of rewrites to achieve their optimization goal.

unclear semantics of macro rules unclear how to define macro rules

How do we optimize programs today?

### **Functional Domain-Specific Compilers**



CGO'18 CGO'17 CASES'16 GPGPU'16 ICFP'15 ...we have developed a **simple automatic search** strategy

...

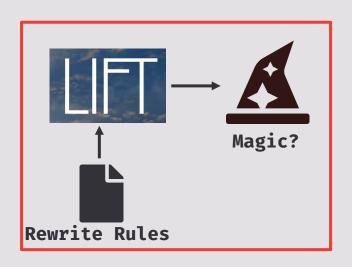
Our current search strategy is *rather basic* and just designed to prove that it is possible to find good implementations automatically.

•••

We *envision replacing this* exploration strategy in the future

How do we optimize programs today?

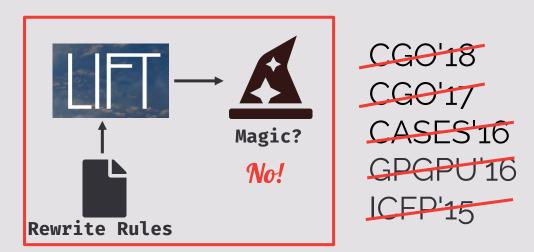
### **Functional Domain-Specific Compilers**



CGO'18 CGO'17 CASES'16 GPGPU'16 ICFP'15

How do we optimize programs today?

### **Functional Domain-Specific Compilers**



The future research is here...

# STRATEGIES

Optimizing Programs like it's 1998 2019

<sup>\*</sup> Elevate's design is inspired by Eelco Visser and others work, e.g. the 1998 ICFP paper "Building program optimizers with rewriting strategies"

# ELEVATE

#### lessons from the past

- A Strategy is a function: Program → Program
- A *Transformation* is the simplest strategy:

```
map(f) \rightarrow join \circ map(map(f)) \circ split(n)
```

• is Defined tests if a strategy is defined for a given program

```
isDefined: Strategy \rightarrow Program \rightarrow Bool
```

• *apply* applies a strategy at a particular location in the program (and might fail)

```
apply: Strategy \rightarrow Location \rightarrow Strategy
```

A *location* and *Traversal* are ADTs:

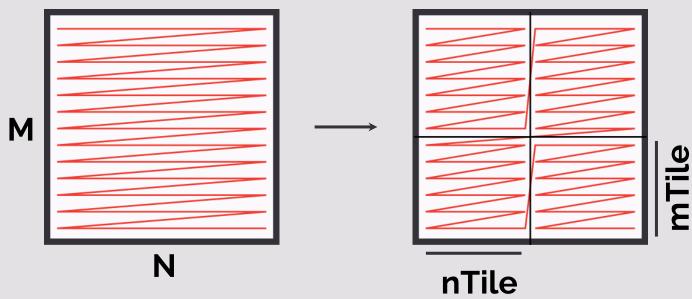
```
data Location = Position(Traversal, Int) | FindFirst(Traversal, Program \rightarrow Bool) data Traversal = BFS | DFS
```

# COMPOSING SRATEGIES

#### defining simple building blocks

```
id: Strategy
id = \lambda p \cdot p
seq: Strategy → Strategy -> Strategy
seq = \lambda f \cdot \lambda s \cdot \lambda p \cdot s (f p)
leftChoice: Strategy → Strategy → Strategy
leftChoice = λf . λs . λp . try (f p) catch (s p)
try: Strategy → Strategy
try = \lambda s. leftChoice s id
repeat: Strategy → Strategy
repeat = \lambda s \cdot try (s ; (repeat s))
normalize: Strategy → Strategy
normalize: λs . repeat(apply s FindFirst(BFS, isDefined s))
```

what is tiling and why is it important for performance



#### Benefits of tiling:

- Exposes more parallelism
- Enables to exploit *locality*

performance improvements of orders of magnitudes

let's define Halide's . tiling in Elevate

```
tiling: Int → Strategy
tiling = λn . λp .
((tileEveryDimension n); // step 1
rewriteNormalForm; // step 2
rearrangeDimensions) p // step 3
```

let's define Halide's . tiling in Elevate

```
Short form for seq

tiling: Int → Strategy
tiling = λn . λp .
((tileEveryDimension r); // step 1
rewriteNormalForm; // step 2
rearrangeDimensions) p // step 3
```

We have *decomposed* the tiling Strategy into three *conceptual steps* that we define also as Strategies

let's define Halide's . tiling in Elevate

### **Elevate:**



# OpenCL:

```
tiling: Int → Strategy
tiling = λn . λp .
((tileEveryDimension n);
rewriteNormalForm;
rearrangeDimensions) p
```



```
for(int i = 0; i < M; i++) {
  for(int j = 0; j < N; j++) {
   out[i][j] = f(in[i][j]);
}}</pre>
```

let's define Halide's . tiling in Elevate

### **Elevate:**

# Lift:

# OpenCL:

#### step 1

```
tileEveryDimension: Int \rightarrow Strategy tileEveryDimension = \lambda n \cdot \lambda p \cdot
```



```
for(int i = 0; i < M; i++) {
  for(int j = 0; j < N; j++) {
   out[i][j] = f(in[i][j]);
}}</pre>
```

```
(split Join size) = *f \rightarrow J \circ **f \circ S(size)
fold: (List L) \rightarrow (P \rightarrow L \rightarrow P) \rightarrow P
```

findAll: (Program → Bool) → Program → (List Location)

let's define Halide's . tiling in Elevate

### **Elevate:**

# Lift:

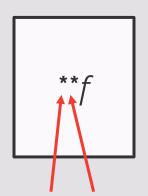
# OpenCL:

#### step 1

```
tileEveryDimension: Int \rightarrow Strategy tileEveryDimension = \lambda n \cdot \lambda p \cdot
```

```
(split Join size) = *f \rightarrow J \circ **f \circ S(size) fold: (List L) \rightarrow (P \rightarrow L \rightarrow P) \rightarrow P \rightarrow P
```

 $findAll: (Program \rightarrow Bool) \rightarrow Program \rightarrow (List \ Location)$ 



```
for(int i = 0; i < M; i++) {
  for(int j = 0; j < N; j++) {
   out[i][j] = f(in[i][j]);
}}</pre>
```

2 locations where split Join rule is defined

let's define Halide's . tiling in Elevate

### **Elevate:**



# OpenCL:

#### step 1

```
tileEveryDimension: Int \rightarrow Strategy tileEveryDimension = \lambda n \cdot \lambda p \cdot
```

```
J。
**(J。**f。S)。
S
```

```
for(int i = 0; i < M; i++) {
  for(int ii = 0; ii < s; ii++) {
    for(int j = 0; j < N; j++) {
     for(int jj = 0; jj < s; jj++) {
        int i_ = i * iTile + ii;
        int j_ = j * jTile + jj;
        out[i_][j_] = f(in[i_][j_]);
}}</pre>
```

```
(splitJoin size) = {}^*f \rightarrow {\bf J} \circ {}^{**}f \circ {\bf S}(size) fold: (List L) \rightarrow (P \rightarrow L \rightarrow P) \rightarrow P
```

findAll: (Program  $\rightarrow$  Bool)  $\rightarrow$  Program  $\rightarrow$  (List Location)

let's define Halide's . tiling in Elevate

### **Elevate:**

### Lift:

# OpenCL:

#### step 1

```
tileEveryDimension: Int \rightarrow Strategy tileEveryDimension = \lambda n \cdot \lambda p \cdot
```

```
J ∘
**(J ∘ **f ∘ S) ∘
S
```

```
for(int i = 0; i < M; i++) {
  for(int ii = 0; ii < s; ii++) {
    for(int j = 0; j < N; j++) {
      for(int jj = 0; jj < s; jj++) {
        int i_ = i * iTile + ii;
        int j_ = j * jTile + jj;
      out[i_][j_] = f(in[i_][j_]);
}
}}</pre>
```

tiling the M & N dimension

```
(split Join size) = *f \rightarrow J \circ **f \circ S(size)
fold: (List L) \rightarrow (P \rightarrow L \rightarrow P) \rightarrow P
```

findAll: (Program  $\rightarrow$  Bool)  $\rightarrow$  Program  $\rightarrow$  (List Location)

let's define Halide's . tiling in Elevate

### **Elevate:**



### OpenCL:

#### step 2

rewriteNormalForm: Strategy rewriteNormalForm = λp . (normalize mapFission) p

mapFission = \*
$$(f \circ g) \rightarrow *f \circ *g$$



```
for(int i = 0; i < M; i++) {
  for(int ii = 0; ii < s; ii++) {
    for(int j = 0; j < N; j++) {
     for(int jj = 0; jj < s; jj++) {
        int i_ = i * iTile + ii;
        int j_ = j * jTile + jj;
        out[i_][j_] = f(in[i_][j_]);
}}</pre>
```

let's define Halide's . tiling in Elevate

### **Elevate:**



# OpenCL:

```
rewriteNormalForm: Strategy
rewriteNormalForm = λp . repeat
(apply mapFission
findFirst (isDefined mapFission))
p
```

```
mapFission = *(f \circ g) \rightarrow *f \circ *g
```

```
J ∘
**(J ∘ **f ∘ S) ∘
S
```

```
for(int i = 0; i < M; i++) {
  for(int ii = 0; ii < s; ii++) {
    for(int j = 0; j < N; j++) {
     for(int jj = 0; jj < s; jj++) {
        int i_ = i * iTile + ii;
        int j_ = j * jTile + jj;
        out[i_][j_] = f(in[i_][j_]);
}}</pre>
```

let's define Halide's . tiling in Elevate

### **Elevate:**



### OpenCL:

```
rewriteNormalForm: Strategy
rewriteNormalForm = λp . repeat
(apply mapFission
findFirst (isDefined mapFission))
p
```

```
mapFission = *(f \circ g) \rightarrow *f \circ *g
```

```
J。**J。
****f。
**S。S
```

```
for(int i = 0; i < M; i++) {
  for(int ii = 0; ii < s; ii++) {
    for(int j = 0; j < N; j++) {
     for(int jj = 0; jj < s; jj++) {
        int i_ = i * iTile + ii;
        int j_ = j * jTile + jj;
        out[i_][j_] = f(in[i_][j_]);
}}
</pre>
```

let's define Halide's . tiling in Elevate

### **Elevate:**



### OpenCL:

```
rearrangeDimensions: Int → Strategy rearrangeDimensions = λd . λp . d match case <2 : p case 2 : (shuffleDimension d) p case _ : (rearrangeDimension (d-1); (shuffleDimension d)) p
```

```
J o ** J o

**** f o

**S o S
```

```
for(int i = 0; i < M; i++) {
  for(int ii = 0; ii < s; ii++) {
    for(int j = 0; j < N; j++) {
     for(int jj = 0; jj < s; jj++) {
        int i_ = i * iTile + ii;
        int j_ = j * jTile + jj;
        out[i_][j_] = f(in[i_][j_]);
}}</pre>
```

let's define Halide's . tiling in Elevate

### **Elevate:**



### OpenCL:

```
rearrangeDimensions: Int → Strategy rearrangeDimensions = λd . λp . d match case <2 : p case 2 : (shuffleDimension d) p case _ : (rearrangeDimension (d-1); (shuffleDimension d)) p
```

```
J∘**J∘
*T∘****f∘*T∘
**S∘S
```

```
for(int i = 0; i < M; i++) {
  for(int ii = 0; ii < s; ii++) {
    for(int j = 0; j < N; j++) {
     for(int jj = 0; jj < s; jj++) {
        int i_ = i * iTile + ii;
        int j_ = j * jTile + jj;
        out[i_][j_] = f(in[i_][j_]);
}}</pre>
```

let's define Halide's . tiling in Elevate

### **Elevate:**

### Lift:

### OpenCL:

#### step 3

```
rearrangeDimensions: Int → Strategy rearrangeDimensions = λd . λp . d match case <2 : p case 2 : (shuffleDimension d) p case _ : (rearrangeDimension (d-1); (shuffleDimension d)) p
```



```
for(int i = 0; i < M; i++) {
  for(int j = 0; j < N; j++) {
   for(int ii = 0; ii < s; ii++) {
    for(int jj = 0; jj < s; jj++) {
      int i_ = i * iTile + ii;
      int j_ = j * jTile + jj;
      out[i_][j_] = f(in[i_][j_]);
}}}</pre>
```

swapped loop-order

Why is defining tiling as a strategy a good idea?

#### ELEVATE

```
tiling: Int → Strategy
tiling = \lambda s \cdot \lambda p.
 ((tileEveryDimension s);
                           // step 1
  rewriteNormalForm; // step 2
   rearrangeDimensions) p // step 3
```

Halide

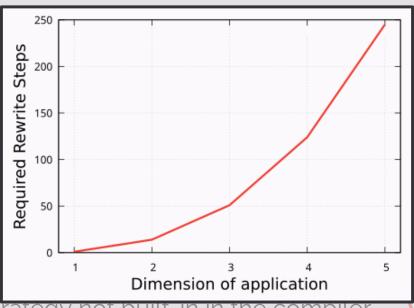
```
out.bound(x, 0, size)
       .bound(y, 0, size)
       .tile(x, y, xi, yi, x_tile * vec_size * ...)
      .split(y, ty, yi, y_unroll)
VS
       .vectorize(xi, vec size)
       .split(xi, xio, xii, warp size)
```

- Strategy not built-in in the compiler
- Works for arbitrary dimensions 2D specific

Why is defining tiling as a strategy a good idea?

### ELEVATE

tiling: Int → Strategy
tiling = λs . λp .
((tileEveryDimension
rewriteNormalForm
rearrangeDimension



#### Halide

```
size)
ze)
yi, x_tile * vec_size * ...)
/i, y_unroll)
/ec_size)
xii, warp_size)
```

Strategy not built-in in the compiler

Strategies are crucial for tiling in higher dimensions

Why is defining tiling as a strategy a good idea?

#### ELEVATE

```
overlapTiling: Int \rightarrow Strategy overlapTiling = \lambda s . \lambda p . ((tileEveryDimension s); // step 1 rewriteNormalForm; // step 2 rearrangeDimensions); // step 3 (try overlapping) p // step 4
```

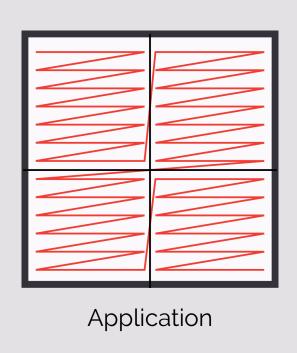
#### Halide

```
vs

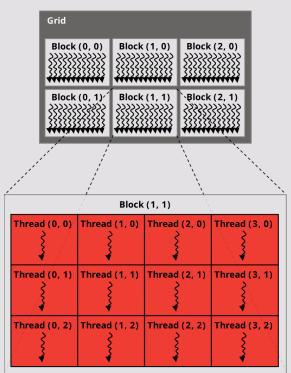
out.bound(x, 0, size)
.bound(y, 0, size)
.tile(x, y, xi, yi, x, le * vec_size * ...)
.split(yi, ty, yi, y nroll)
.vectorize(xi, vec_size)
.split(xi, xio, xii narp_size)
...
```

- Strategy not built-in in the compiler
- Works for arbitrary dimensions
- Easy to extend and reuse existing strategies

#### How to best exploit parallelism in the hardware

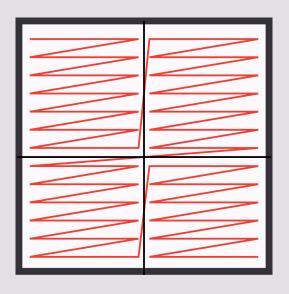






GPU Parallelism Model

#### How to best exploit parallelism in the hardware



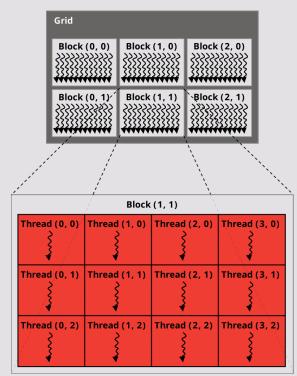
**Application** 





#### Goal:

Minimize parallelism while still fully utilize hardware



GPU Parallelism Model

let's define Halide's .tiling

### **Elevate:**

### Lift:

```
workgroupMapping: Strategy
workgroupMapping = λp . (
(applyOutermost (mapWrg 1));
(applyOutermost (mapWrg 0));
(applyOutermost (mapLcl 1));
(applyOutermost (mapLcl 0));
(normalize mapSeq)) p
```

```
J。**J。*T。
****f。
*T。**S。S
```

```
applyOutermost: Strategy \rightarrow Strategy applyOutermost = \lambdas . \lambdap . apply s (FindFirst DFS (isDefined s)) p
```

let's define Halide's .tiling

### Elevate:

```
workgroupMapping: Strategy
workgroup Mapping = \lambda p. (
 (applyOutermost (mapWrg 1));
 (applyOutermost (mapWrg 0));
 (applyOutermost (mapLcl 1));
 (applyOutermost (mapLcl o));
 (normalize mapSeq)) p
```

```
applyOutermost: Strategy → Strategy
applyOutermost = \lambda s . \lambda p .
apply s (FindFirst DFS (isDefined s)) p
```

# **J** ∘ \*\***J** ∘ \***T** ∘ \*T • \*\*S • S

```
for(int i = 0; i < M; i++) {
 for(int j = 0; j < N; j++) {
 for(int ii = 0; ii < s; ii++) {
   for(int jj = 0; jj < s; jj++) {
         // f
}}}}
```

```
(mapWrg i) = map f \rightarrow mapWrg f
(mapLcl i) = map f \rightarrow mapLcl_i f
```

let's define Halide's .tiling

### **Elevate:**

```
workgroupMapping: Strategy
workgroupMapping = λp . (
(applyOutermost (mapWrg 1));
(applyOutermost (mapWrg 0));
(applyOutermost (mapLcl 1));
(applyOutermost (mapLcl 0));
(normalize mapSeq)) p
```

```
applyOutermost: Strategy \rightarrow Strategy applyOutermost = \lambdas . \lambdap . apply s (FindFirst DFS (isDefined s)) p
```

# Lift:

```
J。**J。*T。
****f。
*T。**S。S
```

```
(mapWrg i) = map f \rightarrow mapWrg_i f

(mapLcl i) = map f \rightarrow mapLcl_i f

f must write to memory
```

#### let's define Halide's .tiling

### **Elevate:**

```
workgroupMapping: Strategy
```

workgroupMapping = λp . (
(applyOutermost (mapWrg 1));
(applyOutermost (mapWrg 0));
(applyOutermost (mapLcl 1));
(applyOutermost (mapLcl 0));
(normalize mapSeq)) p

applyOutermost: Strategy  $\rightarrow$  Strategy applyOutermost =  $\lambda$ s .  $\lambda$ p . apply s (FindFirst DFS (isDefined s)) p

# Lift:

# J o \*\*J o \*T o map(map( map(map f))) o \*T o \*\*S o S

```
(mapWrg i) = map f \rightarrow mapWrg_i f

(mapLcl i) = map f \rightarrow mapLcl_i f
```

let's define Halide's .tiling

### Elevate:

### *Lift:*

```
workgroupMapping: Strategy
workgroup Mapping = \lambda p. (
 (applyOutermost (mapWrg 1));
 (applyOutermost (mapWrg 0));
 (applyOutermost (mapLcl 1));
 (applyOutermost (mapLcl o));
 (normalize mapSeq)) p
```

```
applyOutermost: Strategy → Strategy
applyOutermost = \lambda s . \lambda p .
apply s (FindFirst DFS (isDefined s)) p
```

```
J o **J o *T o
mapWrg<sub>1</sub>(map(
map(map f))) •
*T o **S o S
```

```
int i = get_group_id(1);
 for(int j = 0; j < N; j++) {
  for(int ii = 0; ii < s; ii++) {
   for(int jj = 0; jj < s; jj++) {
}}}
```

```
(mapWrg i) = map f \rightarrow mapWrg f
(mapLcl i) = map f \rightarrow mapLcl_i f
```

let's define Halide's .tiling

### **Elevate:**

### Lift:

```
workgroupMapping: Strategy
workgroupMapping = λp . (
(applyOutermost (mapWrg 1));
(applyOutermost (mapWrg 0));
(applyOutermost (mapLcl 1));
(applyOutermost (mapLcl 0));
(normalize mapSeq)) p
```

```
applyOutermost: Strategy \rightarrow Strategy applyOutermost = \lambdas . \lambdap . apply s (FindFirst DFS (isDefined s)) p
```

```
Jo**Jo*To
mapWrg<sub>1</sub>(mapWrg<sub>0</sub>(
mapLcl<sub>1</sub>(mapLcl<sub>0</sub> f))))
o*To**SoS
```

```
(mapWrg i) = map f \rightarrow mapWrg_i f
(mapLcl i) = map f \rightarrow mapLcl_i f
```

let's define Halide's .tiling

### **Elevate:**

Lift:

### OpenCL:

```
workgroupMapping: Strategy
workgroupMapping = λp . (
(applyOutermost (mapWrg 1));
(applyOutermost (mapWrg 0));
(applyOutermost (mapLcl 1));
(applyOutermost (mapLcl 0));
(normalize mapSeq)) p
```

```
map(map f))
```

```
for(int i = 0; i < M; i++) {
  for(int j = 0; j < N; j++) {
     out[i][j] = f(in[i][j]);
}}</pre>
```

What if we applied this strategy to our non-tiled Lit program?

let's define Halide's .tiling

### Elevate:

```
Lift:
```

### OpenCL:

```
workgroupMapping: Strategy
workgroupMapping = λp.(
(applyOutermost (mapWrg 1));
(applyOutermost (mapWrg 0));
(applyOutermost (mapLcl 1));
(applyOutermost (mapLcl 0));
(normalize mapSeq)) p
```

```
map(map f))
```

```
for(int i = 0; i < M; i++) {
  for(int j = 0; j < N; j++) {
      out[i][j] = f(in[i][j]);
}}</pre>
```

What if we applied this strategy to our non-tiled Lit program?

let's define Halide's .tiling

### **Elevate:**

### Lift:

### OpenCL:

```
globalMapping: Strategy
globalMapping = λp . (
  (applyOutermost (mapGlb 1));
  (applyOutermost (mapGlb 0));
  (normalize mapSeq)) p
```

```
mapGlb₁(
mapGlb₀ f))
```

```
int i = get_global_id(1);
int j = get_global_id(0);
  out[i][j] = f(in[i][j]);
}}
```

What if we applied this strategy to our non-tiled Lit program?

#### exploiting intentional failing of strategies

Short form for leftChoice

mapParallelism: Strategy mapParallelism =  $\lambda p$ . (workgroup3D +> workgroup2D(+>)global2D +> sequential) p

Trying to *exploit all available parallelism* and *gradually fall back* to strategies which make use of less parallelism

Achieves the same goal as Futhark's incremental flattening

### ELEVATE

#### Specifying Compiler Optimizations:

#### • Principled:

One principled way to understand, write and apply compiler optimizations as Strategies

#### • Extensible:

not be fix and built-in. let programmers define abstractions and inject domain & expert knowledge

#### One Language - Many Optimizations:

Controlling different categories of optimizations with the same language

| re-specifi | C        |
|------------|----------|
|            | e-specyu |

Data-Layout

Computation

Memory

Parallelism