To Pack or Not to Pack: A Generalized Packing Analysis and Transformation

Caio S. Rohwedder - Nathan Henderson João P. L. de Carvalho - Yufei Chen - J. Nelson Amaral

Motivation

a_1	a_2	a_3	a_4	a_5	a_6	a_7	a_8	a_9
b_1	b_2	b_3	b_4	b_5	b_6	b_7	b_8	b_9
c_1	c_2	c_3	c_4	c_5	c_6	c_7	c_8	c_9
d_1	d_2	d_3	d_4	d_5	d_6	d_7	d_8	d_9
e_1	e_2	e_3	e_4	e_5	e_6	e_7	e_8	e_9
f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8	f_9
g_1	g_2	g_3	g_4	g_5	g_6	g_7	g_8	g_9
h_1	h_2	h_3	h_4	h_5	h_6	h_7	h_8	h_9
i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	i_9
j_1	j_2	j_3	j_4	j_5	j_6	j_7	j_8	j_9

a_1	a_2	a_3	a_4	a_5	a_6	a_7	a_8	a_9
b_1	b_2	b_3	b_4	b_5	b_6	b_7	b_8	b_9
c_1	c_2	c_3	c_4	C_5	c_6	C7	C_8	<i>C</i> 9
d_1	d_2	d_3	d_4	d_5	d_6	d_7	d_8	d_9
e_1	e_2	e_3	e_4	e_5	e_6	e_7	e_8	e_9
f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8	f_9
g_1	g_2	g_3	g_4	g_5	g_6	g_7	g_8	g_9
h_1	h_2	h_3	h_4	h_5	h_6	h_7	h_8	h_9
i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	i_9
$_{j_1}$	j_2	j_3	j_4	j_5	j_6	j_7	j_8	j_9

-

a_1	a_2	a_3	a_4	a_5	a_6	a_7	a_8	a_9
b_1	b_2	b_3	b_4	b_5	b_6	b_7	b_8	b_9
c_1	C_2	C_3	c_4	c_5	c_6	C7	C_8	<i>C</i> 9
d_1	d_2	d_3	d_4	d_5	d_6	d_7	d_8	d_9
e_1	e_2	e_3	e_4	e_5	e_6	e_7	e_8	e_9
f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8	f_9
g_1	g_2	g_3	g_4	g_5	g_6	g_7	g_8	g_9
h_1	h_2	h_3	h_4	h_5	h_6	h_7	h_8	h_9
i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	i_9
$_{j_1}$	j_2	j_3	j_4	j_5	j_6	j_7	j_8	j_9

_

a_1	a_2	a_3	a_4	a_5	a_6	a_7	a_8	a_9
b_1	b_2	b_3	b_4	b_5	b_6	b_7	b_8	b_9
c_1	C_2	C_3	c_4	C_5	C_6	c_7	c_8	c_9
d_1	d_2	d_3	d_4	d_5	d_6	d_7	d_8	d_9
e_1	e_2	e_3	e_4	e_5	e_6	e_7	e_8	e_9
f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8	f_9
g_1	g_2	g_3	g_4	g_5	g_6	g_7	g_8	g_9
h_1	h_2	h_3	h_4	h_5	h_6	h_7	h_8	h_9
i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	i_9
$_{j_1}$	j_2	j_3	j_4	j_5	j_6	j_7	j_8	j_9

a_1	a_2	a_3	a_4	a_5	a_6	a_7	a_8	a_9
b_1	b_2	b_3	b_4	b_5	b_6	b_7	b_8	b_9
c_1	C_2	C_3	c_4	C_5	C_6	C7	C_8	C9
d_1	d_2	d_3	d_4	d_5	d_6	d_7	d_8	d_9
e_1	e_2	e_3	e_4	e_5	e_6	e_7	e_8	e_9
f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8	f_9
g_1	g_2	g_3	g_4	g_5	g_6	g_7	g_8	g_9
h_1	h_2	h_3	h_4	h_5	h_6	h_7	h_8	h_9
i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	i_9
11	in	ia	i.	1-	in	i-	in	ia

a_1	a_2	a_3	a_4	a_5	a_6	a_7	a_8	a_9
b_1	b_2	b_3	b_4	b_5	b_6	b_7	b_8	b_9
c_1	C_2	C_3	c_4	C_5	C_6	C7	C_8	C9
d_1	d_2	d_3	d_4	d_5	d_6	d_7	d_8	d_9
e_1	e_2	e_3	e_4	e_5	e_6	e_7	e_8	e_9
f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8	f_9
g_1	g_2	g_3	g_4	g_5	g_6	g_7	g_8	g_9
h_1	h_2	h_3	h_4	h_5	h_6	h_7	h_8	h_9
i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	i_9
j_1	12	12	j1	15	16	17	18	ja

a_1	a_2	a_3	a_4	a_5	a_6	a_7	a_8	a_9
b_1	b_2	b_3	b_4	b_5	b_6	b_7	b_8	b_9
c_1	C_2	C_3	c_4	C_5	c_6	C_7	C_8	C9
d_1	d_2	d_3	d_4	d_5	d_6	d_7	d_8	d_9
e_1	e_2	e_3	e_4	e_5	e_6	e_7	e_8	e_9
f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8	f_9
g_1	g_2	g_3	g_4	g_5	g_6	g_7	g_8	g_9
h_1	h_2	h_3	h_4	h_5	h_6	h_7	h_8	h_9
i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	i_9
j_1	j_2	j_3	j_4	j_5	j_6	j_7	j_8	j_9

-





Current state of Loop Optimizers

bondhugula/**pluto**

Pluto: An automatic polyhedral parallelizer and locality optimizer



Stvm

halide/Halide

a language for fast, portable data-parallel computation



Current state of Loop Optimizers

- Tiling
- Unrolling
- Interchanging
- Fusion
- etc.

bondhugula/**pluto**

Pluto: An automatic polyhedral parallelizer and locality optimizer

Stvm

halide/Halide

a language for fast, portable data-parallel computation



MLIR

Packing? No.

Packing? No.

Except for: - GEMM-specific

Packing? No.

Except for:

- GEMM-specific
- Handcrafted
- Autotuned

We propose:

- Compiler-level packing to generic inputs

 Analysis to determine what and where to pack



Packing pros:

- Self-interference misses
 TLB entries
- Vectorization

 e_6



Packing pros:

- Self-interference misses
 TLB entries
- Vectorization



Packing cons:

Copy overhead Naive approach → Slowdown

 e_6

GPAT: Generalized Packing Analysis and Transformation

A Simple Example

Pack tensor **A** targeting loop **ForJ**

Pack tensor **A** targeting loop **ForJ**:

- 1. Insert packing loop
- 2. Substitute A for A'
- 3. Insert unpacking loop

(data-layout change)

Pack tensor **A** targeting loop **ForJ**:

- 1. Insert packing loop
- 2. Substitute A for A'
- 3. Insert unpacking loop

(data-layout change)

for(i=0; i<50; i++) $A' = alloc(1 \times 80 \times 100)$ **for**(m=0; m<80; m++) **for**(n=0; n<100; n++) tmp = load A[m][n][i] store tmp, A'[0][m][n] **for**(j=0; j<60; j++) **for**(k=0; k<80; k++) **for**(1=0; 1<100; 1++) a = **load** A[k][l][i] b = load B[1][k][j]prod = **mul** a, b c = **load** C[i][j] sum = **add** c, prod store sum, C[i][j]

Pack tensor **A** targeting loop **ForJ**:

- 1. Insert packing loop
- 2. Substitute A for A'
- 3. Insert unpacking loop

(data-layout change)

for(i=0; i<50; i++) $A' = alloc(1 \times 80 \times 100)$ **for**(m=0; m<80; m++) **for**(n=0; n<100; n++) tmp = load A[m][n][i] store tmp, A'[0][m][n] **for**(j=0; j<60; j++) **for**(k=0; k<80; k++) **for**(1=0; 1<100; 1++) a = **load** A[k][1][i] b = load B[1][k][j]prod = **mul** a, b c = **load** C[i][j] sum = **add** c, prod store sum, C[i][j]

Pack tensor **A** targeting loop **ForJ**:

- 1. Insert packing loop
- 2. Substitute A for A'
- 3. Insert unpacking loop

(data-layout change)

for(i=0; i<50; i++) $A' = alloc(1 \times 80 \times 100)$ **for**(m=0; m<80; m++) **for**(n=0; n<100; n++) tmp = load A[m][n][i] store tmp, A'[0][m][n] **for**(j=0; j<60; j++) **for**(k=0; k<80; k++) **for**(1=0; 1<100; 1++) a = **load** A'[0][k][1] b = load B[1][k][j]prod = **mul** a, b c = load C[i][j]sum = **add** c, prod store sum, C[i][i]

Pack tensor **A** targeting loop **ForJ**:

- 1. Insert packing loop
- 2. Substitute A for A'
- 3. Insert unpacking loop

(data-layout change)

for(i=0; i<50; i++) $A' = alloc(1 \times 80 \times 100)$ **for**(m=0; m<80; m++) **for**(n=0; n<100; n++) tmp = load A[m][n][i] store tmp, A'[0][m][n] **for**(j=0; j<60; j++) **for**(k=0; k<80; k++) **for**(1=0; 1<100; 1++) a = **load** A'[0][k][1] b = load B[1][k][j]prod = **mul** a, b c = load C[i][j]sum = **add** c, prod store sum, C[i][i]

Pack tensor **A** targeting loop **ForJ**:

- 1. Insert packing loop
- 2. Substitute A for A'
- 3. Insert unpacking loop

(data-layout change)

Packing Analysis

Loop-tensor pair: packing candidate

Loop-tensor pair: packing candidate

Loop-tensor pair: packing candidate

for(i=0; i<50; i++)</pre>

- → **for**(j=0; j<60; j++)
 - for(k=0; k<80; k++)</pre>
 - for(l=0; l<100; l++)</pre>
 - → a = **load** A[k][1][i]
 - b = load B[1][k][j]
 - prod = **mul** a, b
 - c = **load** C[i][j]
 - sum = add c, prod
 - store sum, C[i][j]

Loop-tensor pair: packing candidate


Phase 1 - Data Reuse Filter

Overcome copying overhead with reuse

























Phase 3 - Goal Fulfilment



Phase 3 - Goal Fulfilment



Data-layout change reduces stride at an **innermost** loop?

for(i=0; i<50; i++)
for(j=0; j<60; j++)
for(k=0; k<80; k++)
for(l=0; l<100; l++)
a = load A[k][1][i]</pre>

Data-layout change reduces stride at an **innermost** loop?

for(i=0; i<50; i++)
for(j=0; j<60; j++)
for(k=0; k<80; k++)
for[1=0; 1<100; 1++)
a = load A[k][1][i]
a = load A'[0][k][1]</pre>

Data-layout change reduces stride at an **innermost** loop?

- Cache locality
- Vectorization

Data-layout change reduces stride at an **innermost** loop?

- Cache locality
- Vectorization

What about stride reduction at other loops?

- Less TLB entries

Phase 3 - Goal Fulfilment



Packing reduces TLB entries in a loop below L1 dTLB capacity?







1 TLB entry = 50 elements



1 TLB entry = 50 elements



1 TLB entry = 50 elements



Packing reduces TLB entries in a loop below L1 dTLB capacity?



Packing reduces TLB entries in a loop below L1 dTLB capacity?

- For all **loops affected**
- For all **tensors accessed**
- Given packing and possible data-layout change



Phase 4 - Greedy Selection

Final packing selection:

- 1. Sort candidates based on cost-benefit
- 2. Greedy selection

GPAT in **MLIR** Affine

→ mlir-opt -affine-loop-pack input.mlir



Artifact Available

https://doi.org/10.5281/zenodo.7517506

Evaluation











Packing Choice Evaluation - 2mm



Packing Choice Evaluation - 2mm

$$D_{M \times N} = \beta * D_{M \times N} + T_{M \times N} * C_{N \times N}$$

where
$$T_{M \times N} = \alpha * A_{M \times K} * B_{K \times N}$$


Packing Choice Evaluation - 2mm

- L1 dTLB: miss
- L2 TLB: hit
- (lower is better)



Packing Choice Evaluation - 2mm

- Instruction count
- (lower is better)



Polybench Evaluation

- GPAT compared to other approaches?
- Beyond gemm?
- Robust to prior transformations?



Polybench Evaluation

- GPAT compared to other approaches?
- Beyond gemm?
- Robust to prior transformations?

All polybench benchmarks

Polybench Evaluation

- GPAT compared to other approaches?
- Beyond *gemm*?
- Robust to prior transformations?

Tiling engine:

- Affine tiling

- Speedup over Clang -O3
- Benchmark + Tiling target
- Only benchmarks that were packed



- Speedup over Clang -O3
- Benchmark + Tiling target
- Only benchmarks that were packed







81



Polybench
is designed to
evaluate polly



- Affine tiling is very simplistic



 Even when no tiling "X"



- GPAT improves GEMM-like and GEMM-unalike computations



To Pack or Not to Pack: A Generalized Packing Analysis and Transformation

Caio Salvador Rohwedder csalvado@ualberta.ca

Implemented in MLIR
Evaluated in
Polybench against
Polly and Pluto

More Slides

Candidate (**ForJ, A**):

All accesses of **A** must be invariant
to **j**

for(i=0; i<50; i++) **for**(j=0; j<60; j++) **for**(k=0: k<80: k++) **for**(1=0; 1<100; 1++) a = load A[k][1][i]b = load B[1][k][j]prod = **mul** a, b c = load C[i][j]sum = add c, prod store sum, C[i][j]

Candidate (**ForJ, A**):

 All accesses of A must be invariant to j for(i=0; i<50; i++)</pre>

- → for(j=0; j<60; j++)</p>
 - for(k=0; k<80; k++)</pre>
 - for(l=0; l<100; l++)</pre>
 - → a = **load** A[k][1][i]
 - b = load B[1][k][j]
 - prod = **mul** a, b
 - c = **load** C[i][j]
 - sum = add c, prod
 - store sum, C[i][j]

Candidate (**ForJ, A**):

 All accesses of A must be invariant to j

for(i=0; i<50; i++) → **for**(j=0; j<60; j++) **for**(k=0; k<80; k++) **for**(1=0; 1<100; 1++) \implies a = **load** A[k][1][i] b = load B[1][k][j]prod = **mul** a, b c = load C[i][j]sum = **add** c, prod store sum, C[i][j]

Candidate (L, T):

 All accesses of T must be invariant to the IV of L

Candidate (L, T):

 All accesses of T must be invariant to the IV of L, and any IVs that depend on the IV of L

Ensure A' remains resident:

- Footprint of A'
- Footprint of working set of B and C in one iteration of ForJ

Ensure A' remains resident:

- Footprint of A'
- Footprint of working set of B and C in one iteration of ForJ

Ensure A' remains resident:

- Iter j=0



Ensure A' remains resident:

- Iter j=0



Ensure A' remains resident:

- Iter j=0



Ensure A' remains resident:

- Iter j=0



Ensure A' remains resident:

- Iter j=0



Ensure A' remains resident:

- Iter j=1



Ensure A' remains resident:

- Footprint of A'
- **Twice** the footprint of working set of B and C in one iteration of ForJ

Ensure A' remains resident:

- Footprint of A'
- **Twice** the footprint of working set of B and C in one iteration of ForJ



Ensure A' remains resident:

- Footprint of A'
- **Twice** the footprint of working set of B and C in one iteration of ForJ



Ensure A' remains resident:

- Footprint of A'
- **Twice** the footprint of working set of B and C in one iteration of ForJ



Ensure A' remains resident:

- Footprint of A'
- **Twice** the footprint of working set of B and C in one iteration of ForJ



Ensure A' remains resident:

- Footprint of A'
- **Twice** the footprint of working set of B and C in one iteration of ForJ



Ensure A' remains resident:

- Footprint of A'
- **Twice** the footprint of working set of B and C in one iteration of ForJ


Phase 2 - Cache Residency Filter

Ensure T' remains resident:

- Footprint of T'
- Twice the footprint of working set of all tensors in one iteration of L

1. Sort cost-benefit of candidates:

TLB Improvement Footprint of T' (x2)

1. Sort cost-benefit of candidates:

TLB Improvement Footprint of T' (x2)



1. Sort cost-benefit of candidates:

TLB Improvement Footprint of T' (x2)



- 2. Greedy selection, checking:
 - Redundant packings
 - Candidates benefit in presence of previous selection

- 2. Greedy selection, checking:
 - Redundant packings
 - Candidates benefit in presence of previous selection

```
for(i=0; i<50; i++)
for(j=0; j<60; j++)
for(k=0; k<80; k++)
for(l=0; l<100; l++)
a = load A[k][1][i]
b = load B[1][k][j]
prod = mul a, b
c = load C[i][j]
sum = add c, prod
store sum, C[i][j]</pre>
```

- 2. Greedy selection, checking:
 - Redundant packings
 - Candidates benefit in presence of previous selection

Goals



Packing Choice Evaluation - *gemm*



Packing Choice Evaluation - *gemm* **BLIS**

gemm interchanged
 to BLIS loop ordering



- Speedup over Clang O3
- Polly



gemm



mm



3mm

