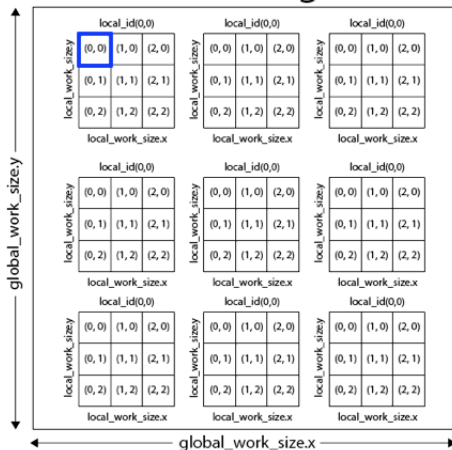# Input Space Splitting for OpenCL

Simon Moll, Johannes Doerfert, Sebastian Hack
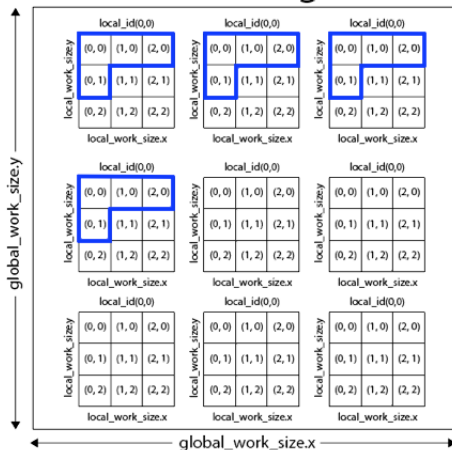
Saarbrücken Graduate School of Computer Science
Saarland University
Saarbrücken, Germany

October 29, 2015

SAARLAND
UNIVERSITY

GRADUATE SCHOOL OF
COMPUTER SCIENCE

# NDRange

SAARLAND
UNIVERSITY
COMPUTER SCIENCE



NDRange

# Vectorization (SIMD)

Perform the same operation for multiple vector lanes simultaneously.

# Vectorization (SIMD)

Perform the same operation for multiple vector lanes simultaneously.

### Vector Patterns

| | | |
|---|---|---|
| Consecutive: | contiguous entries | $< i, i + 1, i + 2, i + 3 >$ |
| Uniform: | single entry | $<i,i,i,i> \rightarrow i$ |
| Divergent: | unrelated entries | $< i, j, 7, - >$ |

```
for (i = 0; i < 16; i++)      for (i = 0; i < 16; i += 2)
  O[i] = I[i] + 2;              O[i] = I[i] + 1;
```

# Diverging Control Flow



| Thread | Trace |
|--------|-------|
| 1 | a b c e f |
| 2 | a b d e f |
| 3 | a b c e b c e f |
| 4 | a b c e b d e f |

- Different threads execute different code paths

# Diverging Control Flow
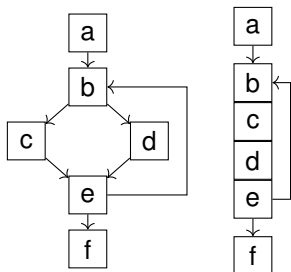


| Thread | Trace |
|--------|-------|
| 1 | a b c d e b c d e f |
| 2 | a b c d e b c d e f |
| 3 | a b c d e b c d e f |
| 4 | a b c d e b c d e f |

- Different threads execute different code paths
- Execute everything, mask out results of inactive threads (using predication, blending)
- Control flow to data flow conversion on ASTs [Allen & Kennedy '83]
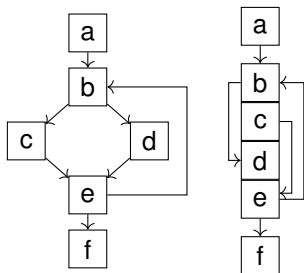- Whole-Function Vectorization on SSA CFGs [K & H '11]

# Non-Divergent Control Flow

- Idea: optimize cases where threads do **not** diverge



| Thread | Trace |
|--------|-------|
| 1 | a b c e b d e f |
| 2 | a b c e b d e f |
| 3 | a b c e b d e f |
| 4 | a b c e b d e f |

# Non-Divergent Control Flow
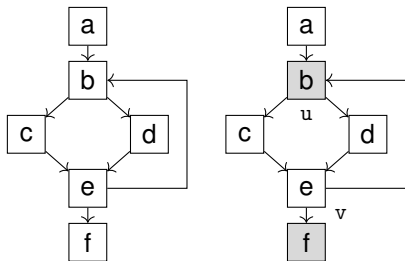
- Idea: optimize cases where threads do not diverge



| Thread | Trace |
|--------|-------------------|
| 1 | a b c e b d e f |
| 2 | a b c e b d e f |
| 3 | a b c e b d e f |
| 4 | a b c e b d e f |

- Option 1: Insert dynamic predicate-tests & branches to skip paths
  - ▸ "Branch on superword condition code" (BOSCC) [Shin et al. PACT'07]
  - ▸ Additional overhead for dynamic test
  - ▸ Does not help against increased register pressure

# Non-Divergent Control Flow
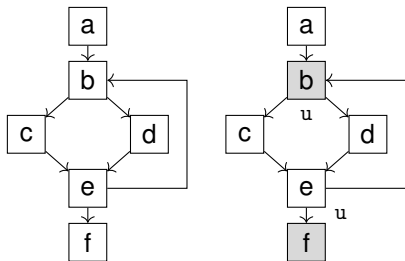
- Idea: optimize cases where threads do not diverge

| Thread | Trace |
|--------|-------|
| 1 | a b c e b d e f |
| 2 | a b c e b d e f |
| 3 | a b c e b d e f |
| 4 | a b c e b d e f |

- Option 2: Statically prove non-divergence of certain blocks
  - Non-divergent blocks can be excluded from linearization
  - Less executed code, less register pressure
  - More conservative than dynamic test

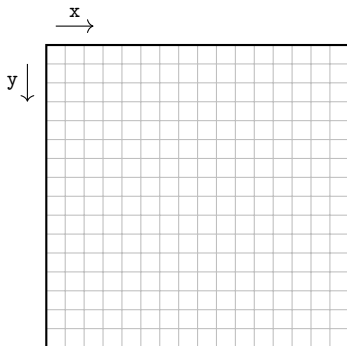# Non-Divergent Control Flow

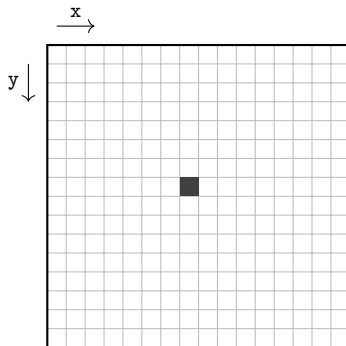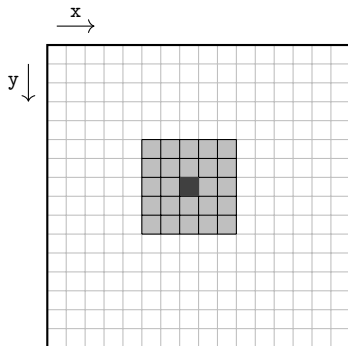- Idea: optimize cases where threads do **not** diverge



| Thread | Trace |
|--------|-------|
| 1 | a b c e f |
| 2 | a b c e f |
| 3 | a b c e b d e f |
| 4 | a b c e b d e f |
| 5 | a b c e b d e f |
| 6 | a b c e b d e f |

- Option 3: Statically split non-divergence inputs
  - ► Code versions with improved divergence properties
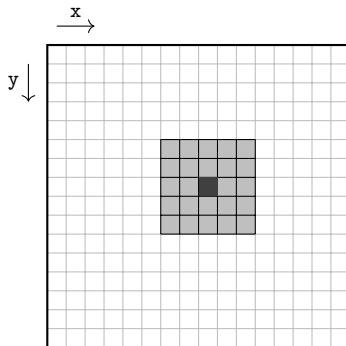  - ► Orthogonal to both other options ⟹ combination possible

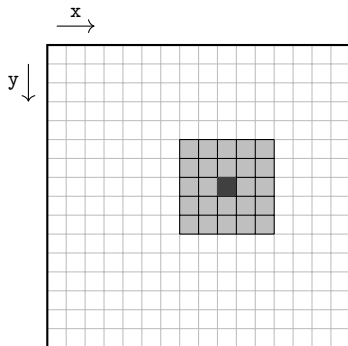# 2D Convolution

# 2D Convolution
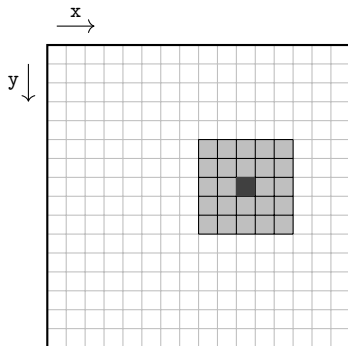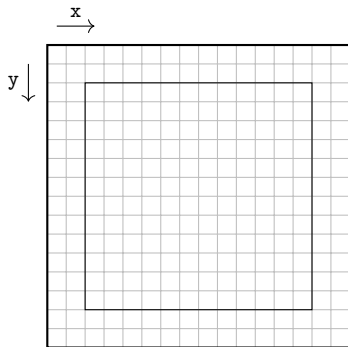
# 2D Convolution

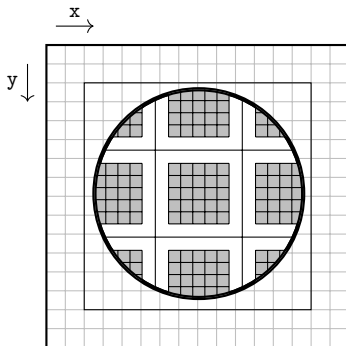# 2D Convolution

# 2D Convolution

```
int left   = x - 2;
int right  = x + 2;
int top    = y - 2;
int bottom = y + 2;

int sum = 0;
for (int i = left; i <= right; ++i)
  for (int j = top; j <= bottom; ++j)
    sum += input[j][i] * mask[j - top][i - left];
output[y][x] = sum;
```

# 2D Convolution

```
auto left   = x - 2;
auto right  = x + 2;
int top     = y - 2;
int bottom  = y + 2;

int sum = 0;
for (auto i = left; i <= right; ++i)
  for (int j = top; j <= bottom; ++j)
    sum += input[j][i] * mask[j - top][i - left];
output[y][x] = sum;
```
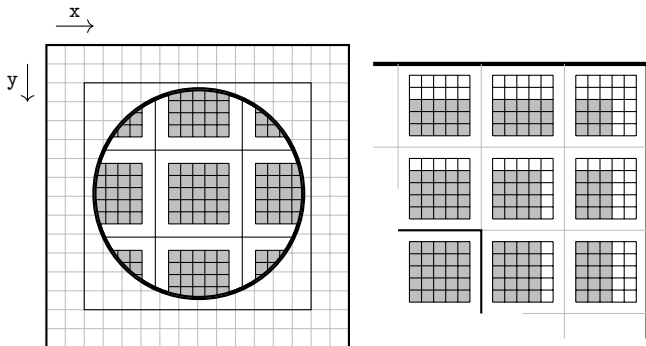
# 2D Convolution

# 2D Convolution
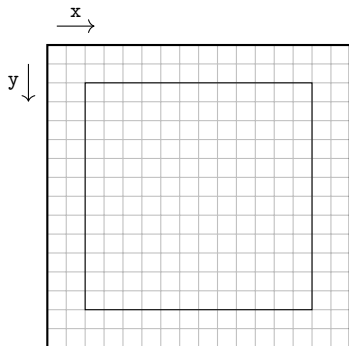


x →

y ↓

# 2D Convolution

```
int left   = MAX(0, x - 2);
int right  = MIN(width - 1, x + 2);
int top    = MAX(0, y - 2);
int bottom = MIN(height - 1, y + 2);

int sum = 0;
for (int i = left; i <= right; ++i)
  for (int j = top; j <= bottom; ++j)
    sum += input[j][i] * mask[j - (y - 2)][i - (x - 2)];
output[y][x] = sum;
```

# 2D Convolution

```
auto left  = MAX(0, x - 2);
auto right = MIN(width - 1, x + 2);
int top    = MAX(0, y - 2);
int bottom = MIN(height - 1, y + 2);

int sum = 0;
for (auto i = left; i <= right; ++i)
  for (int j = top; j <= bottom; ++j)
    sum += input[j][i] * mask[j - (y - 2)][i - (x - 2)];
output[y][x] = sum;
```
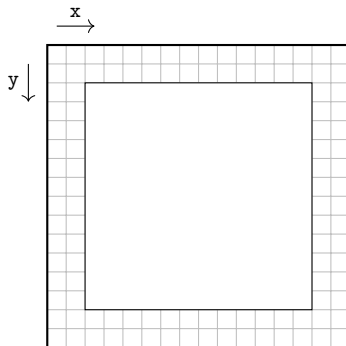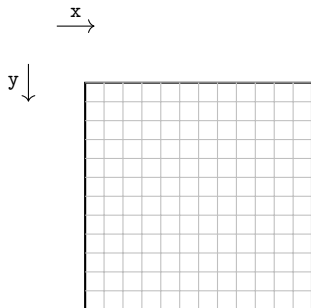
SAARLAND
UNIVERSITY
COMPUTER SCIENCE

```
S:   A[i][j] = /* ... */;
     if (j <= i)
P:     A[i][j]+= A[j][i];
```

```
for (int i = 0; i <= N; i++)
  for (int j = 0; j <= N; j++) {
S:  A[i][j] = /* ... */;
    if (j <= i)
P:    A[i][j]+= A[j][i];
  }
```

```
for (int i = 0; i <= N; i++)
  for (int j = 0; j <= N; j++) {
S:  A[i][j] = /* ... */;
    if (j <= i)
P:    A[i][j]+= A[j][i];
  }
```



$$\mathcal{I}_S = \{(S, (i, j)) \mid 0 \le i \le N \wedge 0 \le j \le N\}$$

# The Polyhedral Model

```
for (int i = 0; i <= N; i++)
  for (int j = 0; j <= N; j++) {
S:  A[i][j] = /* ... */;
    if (j <= i)
P:    A[i][j]+= A[j][i];
  }
```



$$\mathcal{I}_S = \{(S, (i, j)) \mid 0 \le i \le N \land 0 \le j \le N\}$$

$$\mathcal{I}_P = \{(P, (i, j)) \mid 0 \le i \le N \land 0 \le j \le i\}$$

```
for (int i = 0; i <= N; i++)
  for (int j = 0; j <= N; j++) {
S:  A[i][j] = /* ... */;
    if (j <= i)
P:    A[i][j]+= A[j][i];
  }
```
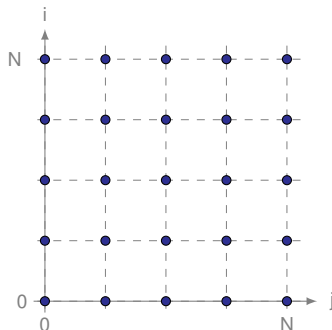


$$\mathcal{F}_S = \{(S, (i, j)) \rightarrow (i, j)\}$$

# The Polyhedral Model
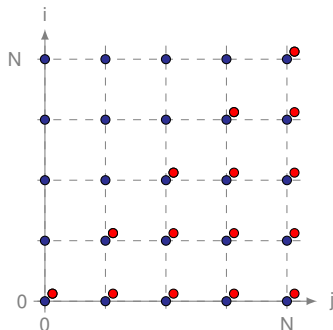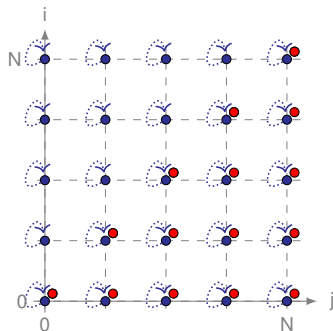
```
for (int i = 0; i <= N; i++)
  for (int j = 0; j <= N; j++) {
S:  A[i][j] = /* ... */;
    if (j <= i)
P:    A[i][j]+= A[j][i];
  }
```



$$\mathcal{F}_S = \{(S, (i, j)) \rightarrow (i, j)\}$$

$$\mathcal{F}_{P_1} = \{(P, (i, j)) \rightarrow (i, j)\} \quad \mathcal{F}_{P_2} = \{(P, (i, j)) \rightarrow (j, i)\}$$

# Splitting Predicates

Full Tile Predicate

```
for (int i = 0; i <= N; i++)
  for (int j = 0; j <= N; j += 8) {
S:  A[i][j] = /* ... */;
    if (j <= i)
P:    A[i][j]+= A[j][i];
  }
```

# Splitting Predicates

Full Tile Predicate

```
for (int i = 0; i <= N; i++)
  for (int j = 0; j <= N; j += 8) {
S:  A[i][j] = /* ... */;
    if (j <= i)
P:    A[i][j]+= A[j][i];
  }
```

$$\mathsf{Full_S} = \{(S, (i,j)) \mid (j - (j \bmod 8)) + 7 \leq N\}$$

# Splitting Predicates

Full Tile Predicate

```
for (int i = 0; i <= N; i++)
  for (int j = 0; j <= N; j += 8) {
S:  A[i][j] = /* ... */;
    if (j <= i)
P:    A[i][j]+= A[j][i];
  }
```

$$\text{Full}_S = \{(S,(i,j)) \mid (j - (j \bmod 8)) + 7 \leq N\}$$
$$\text{Full}_P = \{(P,(i,j)) \mid (j - (j \bmod 8)) + 7 \leq \min(i, N)\}$$

# Splitting Predicates

Uniform Access Predicate

```
for (int i = 0; i <= N; i++)
  for (int j = 0; j <= N; j += 8) {
S:  A[i][j] = /* ... */;
    if (j <= i)
P:    A[i][j]+= A[j][i];
  }
```

$$\text{Uni}_{\mathcal{F}_S} = \{(S, (i,j)) \mid \mathcal{F}_S(i, j+1) = \mathcal{F}_S(i,j)\}$$

# Splitting Predicates

Uniform Access Predicate

```
for (int i = 0; i <= N; i++)
  for (int j = 0; j <= N; j += 8) {
S:  A[i][j] = /* ... */;
    if (j <= i)
P:    A[i][j]+= A[j][i];
  }
```

$$\mathsf{Uni}_{\mathcal{F}_S} = \{(S, (i,j)) \mid \mathcal{F}_S\,(i, j+1) = \mathcal{F}_S\,(i,j)\}$$

$$\mathsf{Uni}_{\mathcal{F}_{P_1}} = \{(P, (i,j)) \mid \mathcal{F}_{P_1}(i, j+1) = \mathcal{F}_{P_1}(i,j)\}$$

$$\mathsf{Uni}_{\mathcal{F}_{P_2}} = \{(P, (i,j)) \mid \mathcal{F}_{P_2}(i, j+1) = \mathcal{F}_{P_2}(i,j)\}$$

# Splitting Predicates

Uniform Access Predicate

```
for (int i = 0; i <= N; i++)
  for (int j = 0; j <= N; j += 8) {
S:  A[i][j] = /* ... */;
    if (j <= i)
P:    A[i][j] += A[j][i];
  }
```

$$\text{Uni}_{\mathcal{F}_S} = \{(S, (i,j)) \mid \mathcal{F}_S(i,j+1) = \mathcal{F}_S(i,j)\} = \{\}$$

$$\text{Uni}_{\mathcal{F}_{P_1}} = \{(P, (i,j)) \mid \mathcal{F}_{P_1}(i,j+1) = \mathcal{F}_{P_1}(i,j)\} = \{\}$$

$$\text{Uni}_{\mathcal{F}_{P_2}} = \{(P, (i,j)) \mid \mathcal{F}_{P_2}(i,j+1) = \mathcal{F}_{P_2}(i,j)\} = \{\}$$

# Splitting Predicates

Consecutive Access Predicate

```
for (int i = 0; i <= N; i++)
  for (int j = 0; j <= N; j += 8) {
S:  A[i][j] = /* ... */;
    if (j <= i)
P:    A[i][j]+= A[j][i];
  }
```

$$\text{Cons}_{\mathcal{F}_S} = \{(S, (i,j)) \mid \mathcal{F}_S\,(i, j+1) = \mathcal{F}_S\,(i, j) + 1\}$$

$$\text{Cons}_{\mathcal{F}_{P_1}} = \{(P, (i,j)) \mid \mathcal{F}_{P_1}(i, j+1) = \mathcal{F}_{P_1}(i, j) + 1\}$$

$$\text{Cons}_{\mathcal{F}_{P_2}} = \{(P, (i,j)) \mid \mathcal{F}_{P_2}(i, j+1) = \mathcal{F}_{P_2}(i, j) + 1\}$$

# Splitting Predicates

Consecutive Access Predicate

```
for (int i = 0; i <= N; i++)
  for (int j = 0; j <= N; j += 8) {
S:  A[i][j] = /* ... */;
    if (j <= i)
P:    A[i][j]+= A[j][i];
  }
```

$$\text{Cons}_{\mathcal{F}_S} = \{(S, (i, j)) \mid \mathcal{F}_S\,(i, j+1) = \mathcal{F}_S\,(i, j) + 1\} = \mathcal{I}_S$$

$$\text{Cons}_{\mathcal{F}_{P_1}} = \{(P, (i, j)) \mid \mathcal{F}_{P_1}(i, j+1) = \mathcal{F}_{P_1}(i, j) + 1\} = \mathcal{I}_P$$

$$\text{Cons}_{\mathcal{F}_{P_2}} = \{(P, (i, j)) \mid \mathcal{F}_{P_2}(i, j+1) = \mathcal{F}_{P_2}(i, j) + 1\} = \{\}$$

```
for (int i = 0; i <= N; i++)
  for (int j = 0; j <= N; j += 8)
    if (i <= NumParticles) {
S:    ...
    }
```

# CFG simplification

Hoisting conditionals

```
for (int i = 0; i <= N; i++)
  for (int j = 0; j <= N; j += 8)
    if (i <= NumParticles) {
S:    ...
    }



for (int i = 0; i <= NumParticles; i++)
  for (int j = 0; j <= i; j += 8)
S:  ...
```

# CFG simplification

Hoisting conditionals

```
for (int i = 0; i <= N; i++)
  for (int j = 0; j <= N; j += 8)
    if (reverse) {
S:    ...
    } else {
P:    ...
    }
```
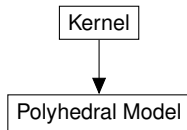
Hoisting conditionals

```
for (int i = 0; i <= N; i++)
  for (int j = 0; j <= N; j += 8)
    if (reverse) {
S:    ...
    } else {
P:    ...
    }

if (reverse) {
  for (int i = 0; i <= N; i++)
    for (int j = 0; j <= i; j += 8)
  S: ...
} else {
  for (int i = 0; i <= N; i++)
    for (int j = 0; j <= i; j += 8)
  P: ...
}
```
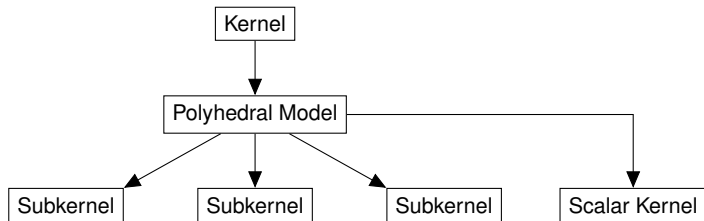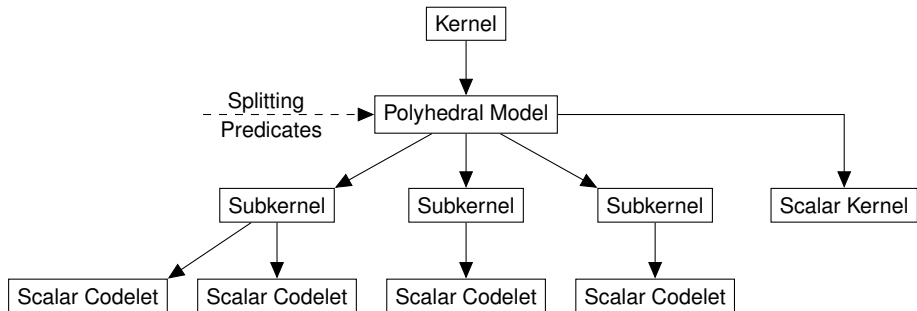
# Predicate Based Domain Splitting

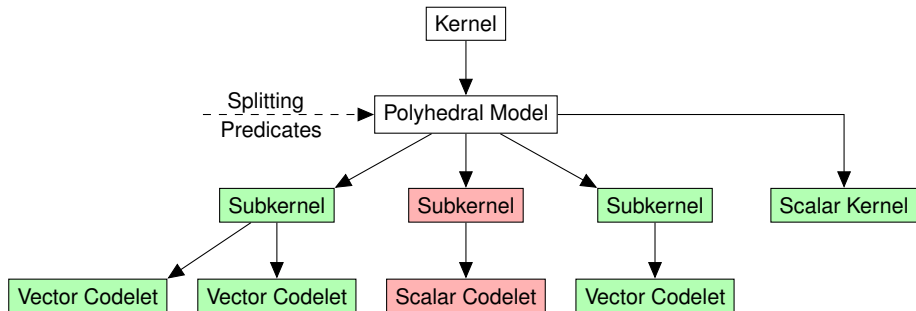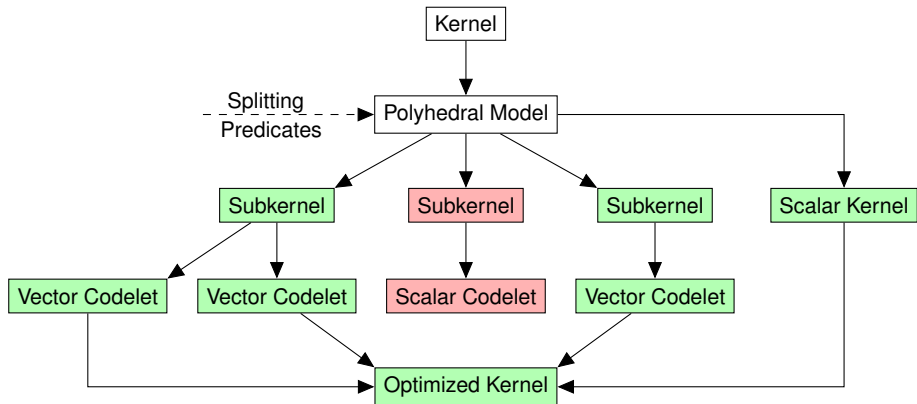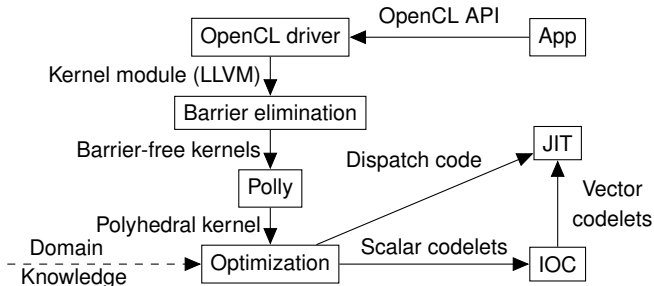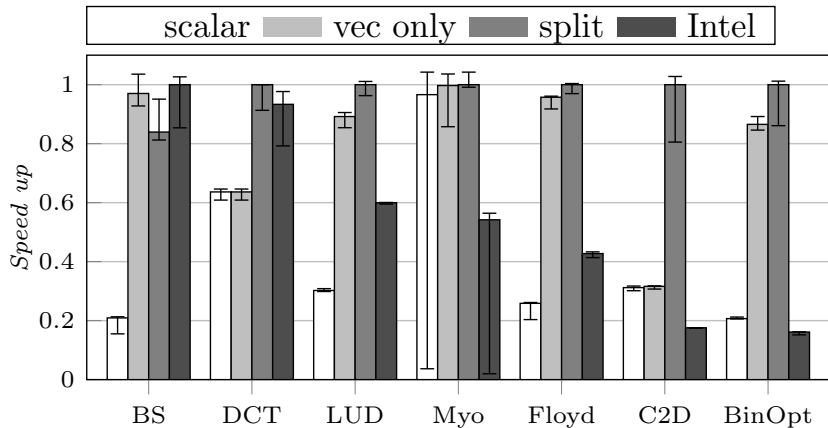Kernel

# Predicate Based Domain Splitting

```
┌─────────┐
│ Kernel  │
└─────────┘
     │
     ▼
┌──────────────────┐
│ Polyhedral Model │
└──────────────────┘
```

# Predicate Based Domain Splitting

```
                    ┌──────────┐
                    │  Kernel  │
                    └──────────┘
                          │
                          ▼
                ┌──────────────────┐
                │ Polyhedral Model │
                └──────────────────┘
                   │      │      │
          ┌────────┘      │      └────────┐
          ▼               ▼               ▼
   ┌───────────┐   ┌───────────┐   ┌───────────┐
   │ Subkernel │   │ Subkernel │   │ Subkernel │
   └───────────┘   └───────────┘   └───────────┘
```

# Predicate Based Domain Splitting

# Predicate Based Domain Splitting
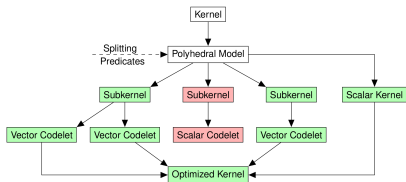
# Predicate Based Domain Splitting
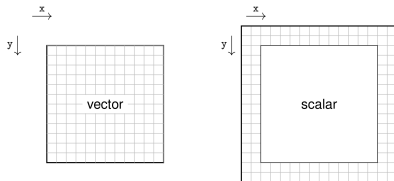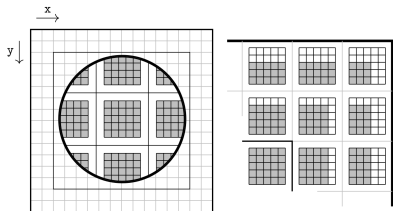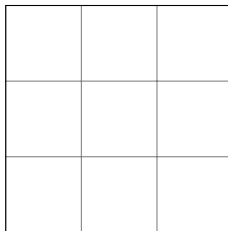
# Evaluation

Performance

- Model synchronization the Polyhedral Model.
- Apply polyhedral optimizations (scheduling).
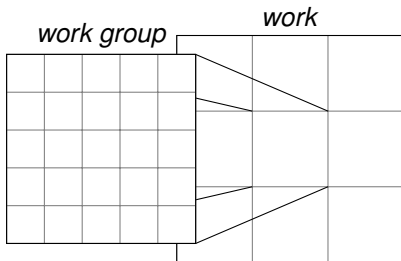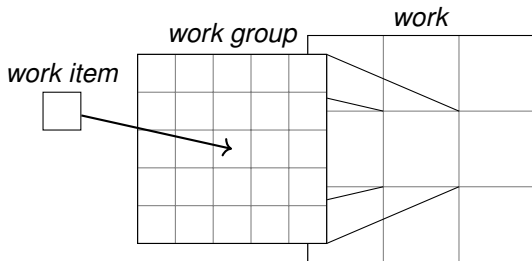- Improve the representation of non-affine parts.

# Conclusion

# OpenCL Programming Model

*work*

*work group*   *work*

# Codelet Score

$$Score_n(k) := \begin{cases} \Sigma_{Q \in k, \mathcal{F} \in \mathbb{F}_Q} \\ w_{cons}\|Box(Cons_{\mathcal{F}}(d_k))\| & \text{if } n \geq w \\ +w_{uni}\|Box(Uni_{\mathcal{F}}(d_k))\| \\ 0 & otw. \end{cases}$$

$$\mathcal{I}_k^{\mathcal{C}} := \bigcap_{Q \in k} \bigcap_{\substack{\mathcal{F} \in \mathbb{F}_Q, st \\ Cons_{\mathcal{F}}(d_k) \neq \varnothing}} Cons_{\mathcal{F}}(d_k) \,.$$

SAARLAND
UNIVERSITY
COMPUTER SCIENCE