

Matrix Multiplication Beyond Auto-Tuning: Rewrite Based GPU Code Generation



Michel Steuwer



Toomas Rimmelg



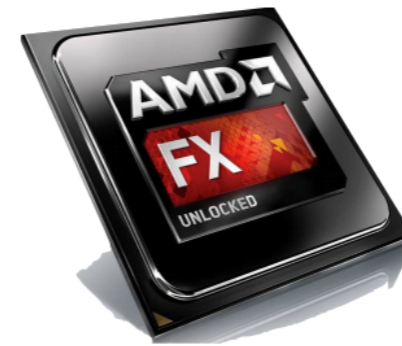
Christophe Dubach



THE UNIVERSITY
of EDINBURGH

Motivation

- Tuning parallel programs is complex
- The diversity of heterogeneous accelerators makes it even harder
- Auto-tuning has been successfully applied to ease this
- **Problem of Performance-Portability:**
 - Parametric kernels are complex
 - Still falls short on new architectures!



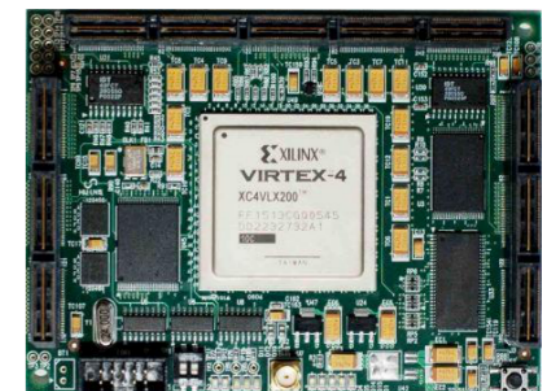
CPU



GPU



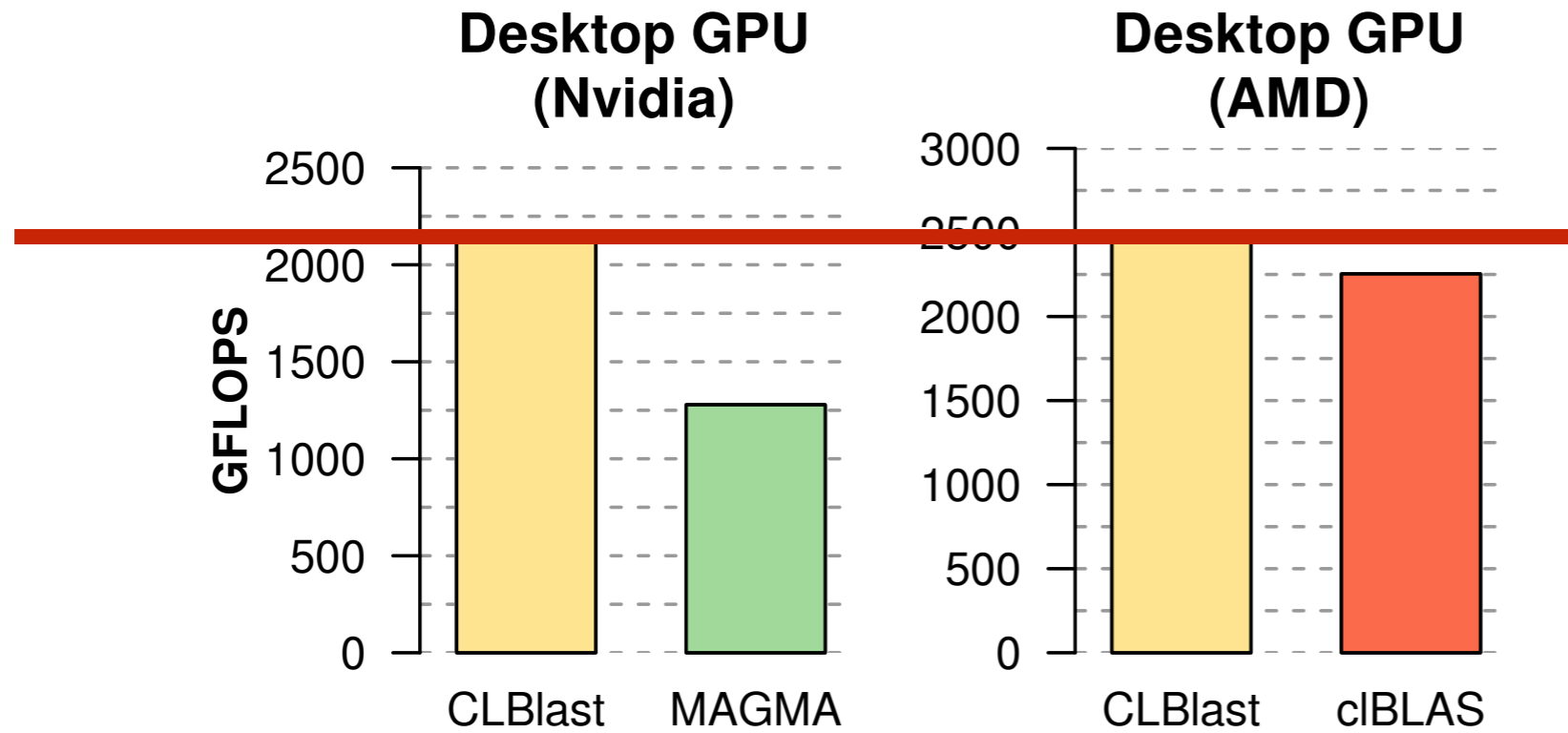
Accelerator



FPGA

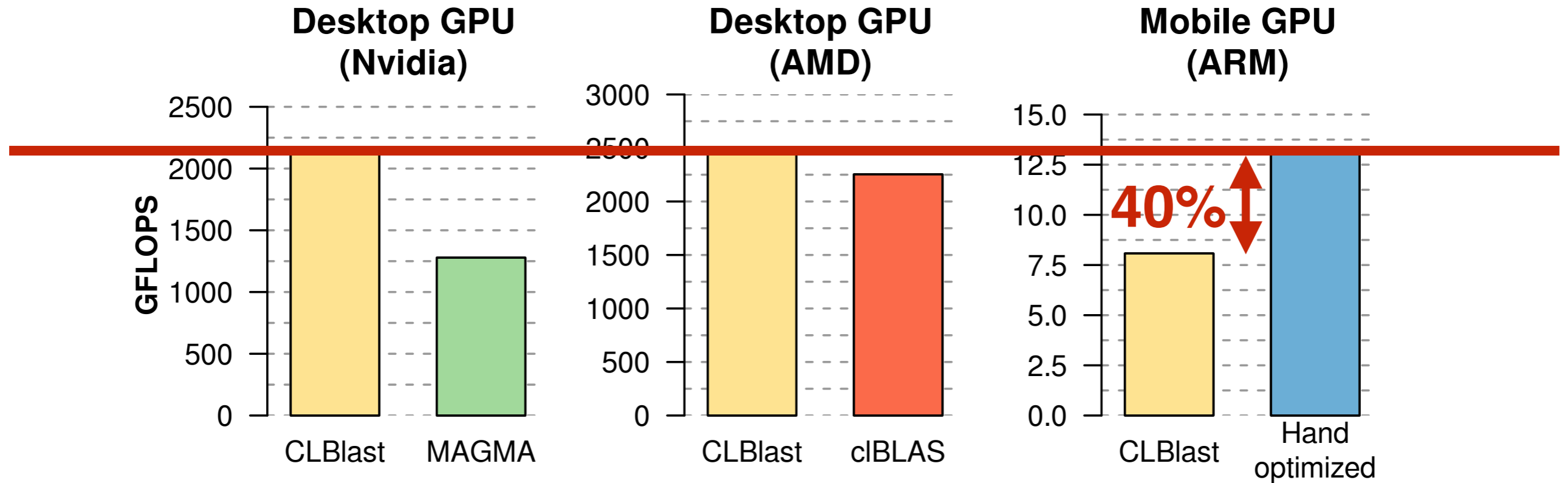
The Problem of Performance-Portability

Matrix-Multiplication



The Problem of Performance-Portability

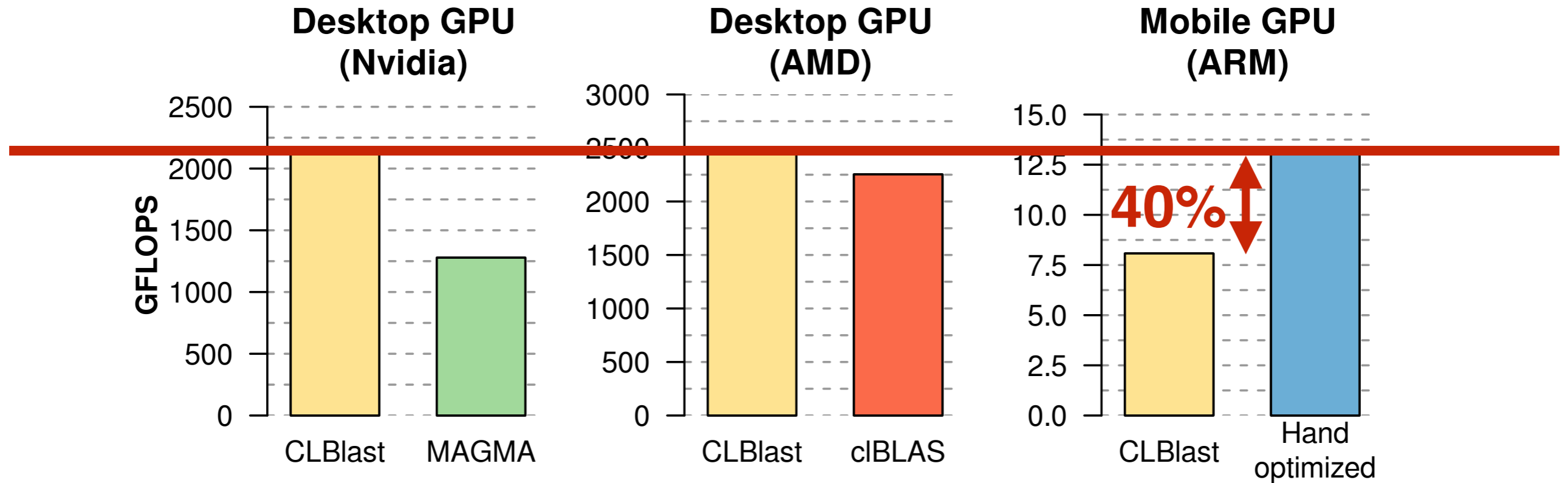
Matrix-Multiplication



- Auto-tuning fails to deliver on a significantly different architecture

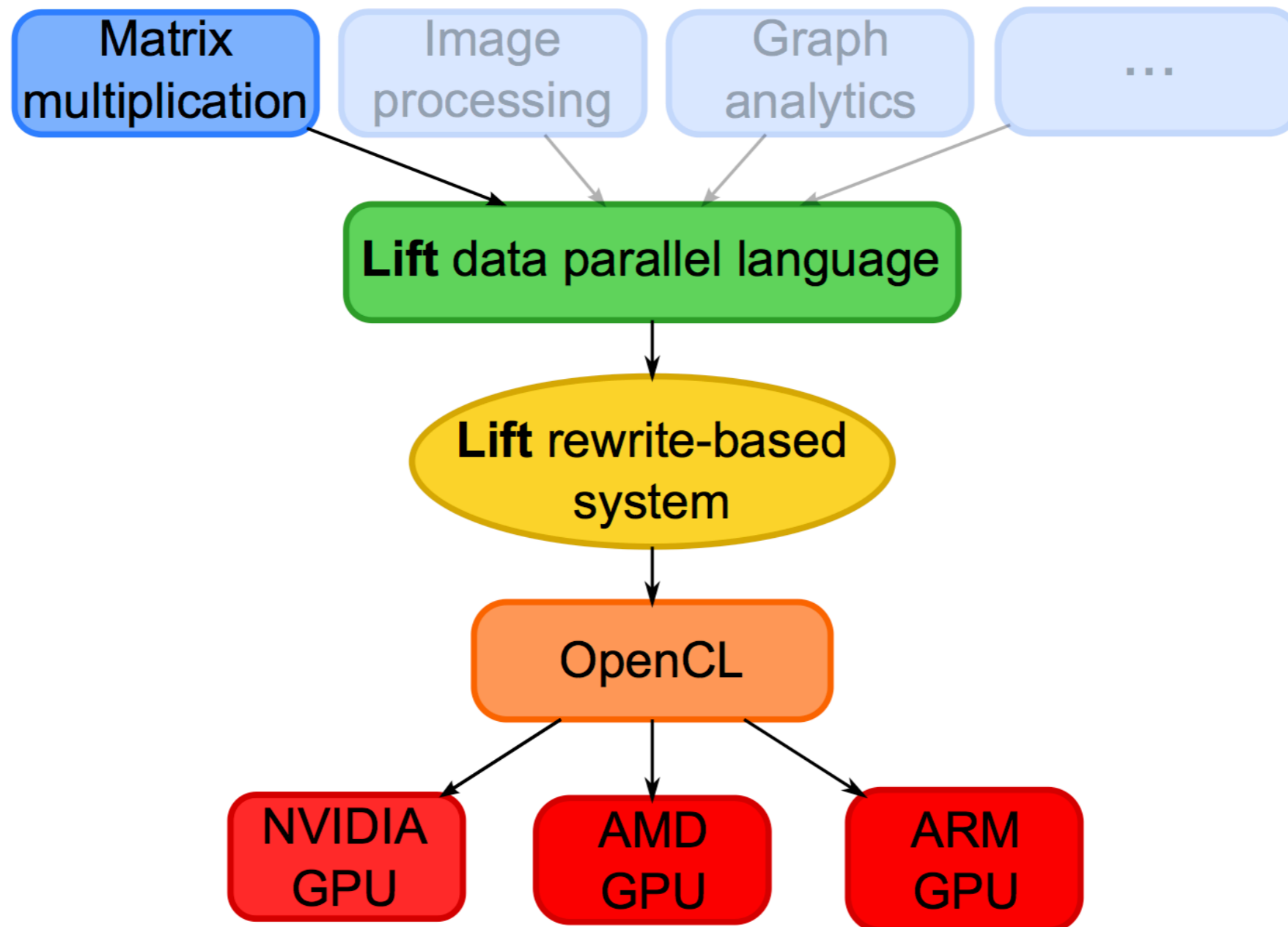
The Problem of Performance-Portability

Matrix-Multiplication



- Auto-tuning fails to deliver on a significantly different architecture
- **Limitation of Auto-tuning:**
Fixed set of parameters and optimisations

Our Approach to Performance Portability



Lift Data Parallel Language

$$\text{map}(f, \boxed{x_1 \ x_2 \ \dots \ x_n}) = \boxed{f(x_1) \ f(x_2) \ \dots \ f(x_n)}$$

$$\text{reduce}(z, f, \boxed{x_1 \ x_2 \ \dots \ x_n}) = \boxed{f(\dots (f(f(z, x_1), x_2) \dots), x_n)}$$

$$\text{zip}(\boxed{x_1 \ x_2 \ \dots \ x_n}, \boxed{y_1 \ y_2 \ \dots \ y_n}) = \boxed{(x_1, y_1) \ (x_2, y_2) \ \dots \ (x_n, y_n)}$$

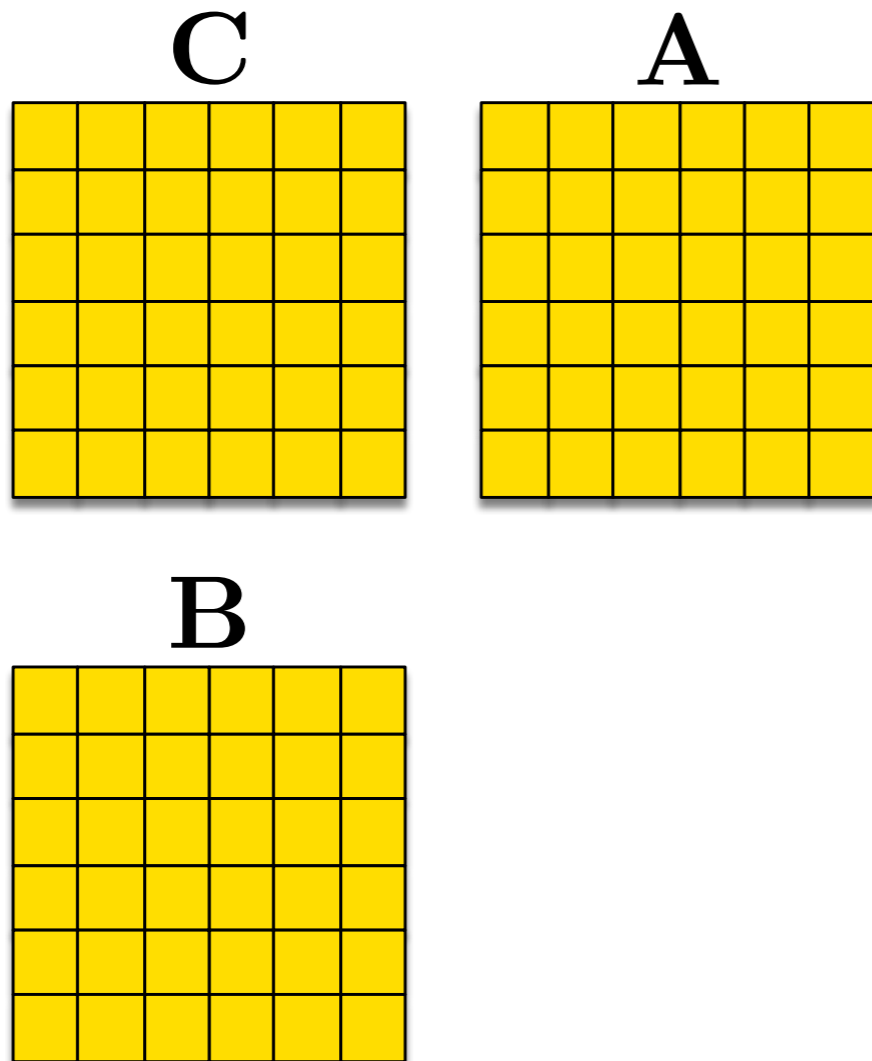
$$\text{id}(\boxed{x_1 \ x_2 \ \dots \ x_n}) = \boxed{x_1 \ x_2 \ \dots \ x_n}$$

$$\text{split}^m(\boxed{x_1 \ x_2 \ \dots \ \dots \ \dots \ \dots \ \dots \ \dots \ \dots \ x_n}) = \boxed{\boxed{x_1 \ x_2 \ \dots \ \dots} \ \boxed{\dots \ \dots \ \dots \ \dots} \ \dots \ \boxed{\dots \ \dots \ \dots \ x_n}}$$

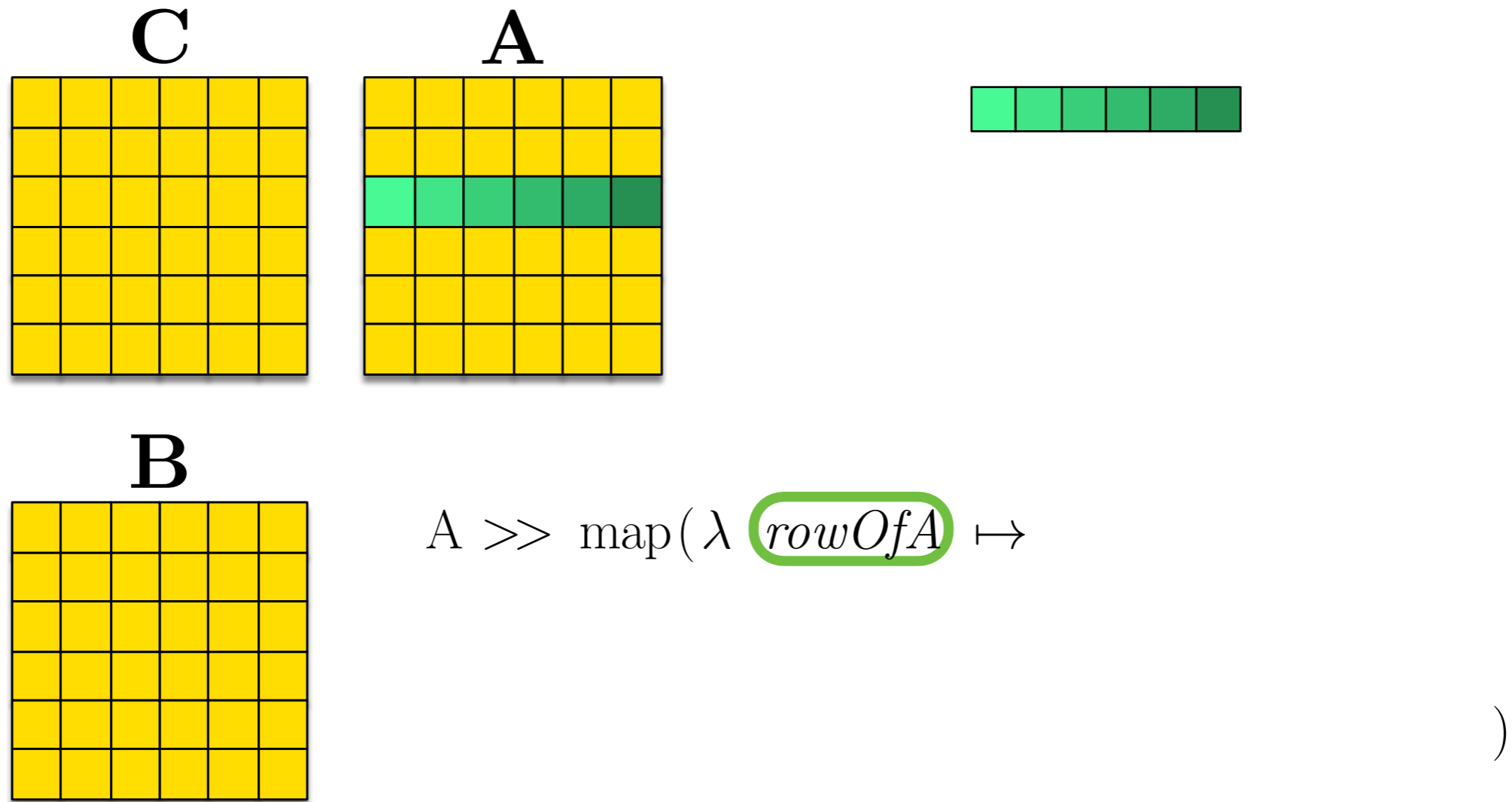
$\underbrace{\hspace{10em}}_m$

$$\text{join}(\boxed{\boxed{x_1 \ x_2 \ \dots \ \dots} \ \boxed{\dots \ \dots \ \dots \ \dots} \ \dots \ \boxed{\dots \ \dots \ \dots \ x_n}}) = \boxed{x_1 \ x_2 \ \dots \ \dots \ \dots \ \dots \ \dots \ \dots \ \dots \ x_n}$$

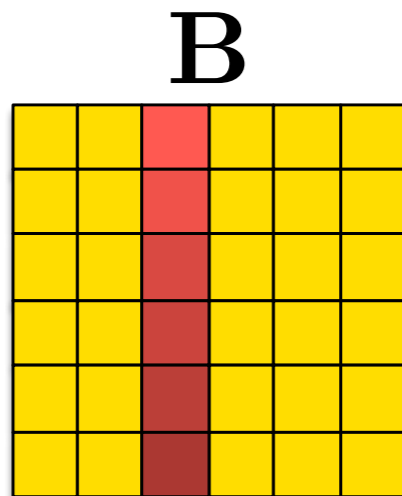
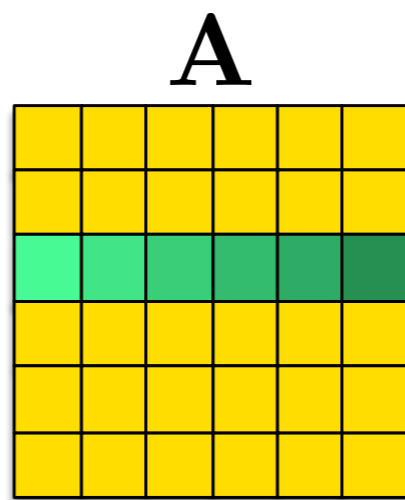
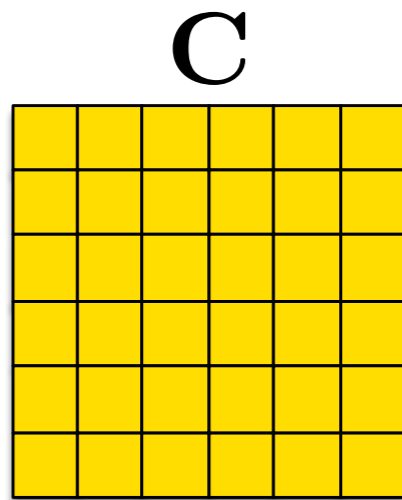
Matrix Multiplication Expressed Functionally



Matrix Multiplication Expressed Functionally



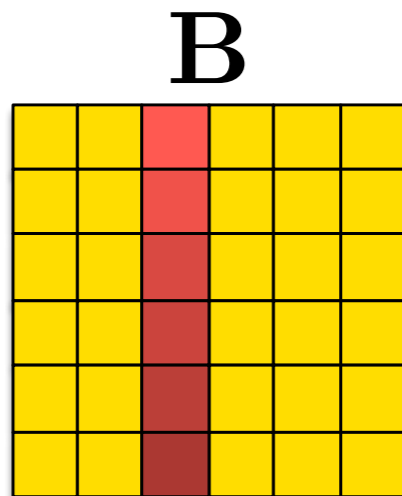
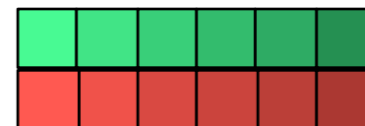
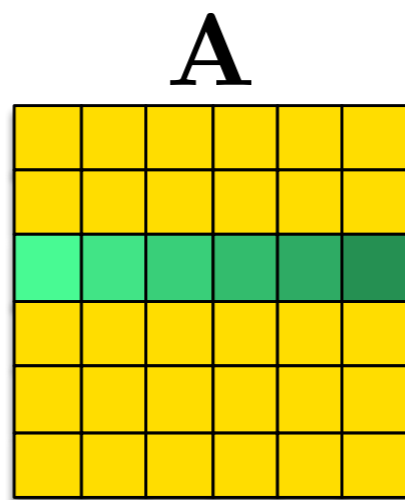
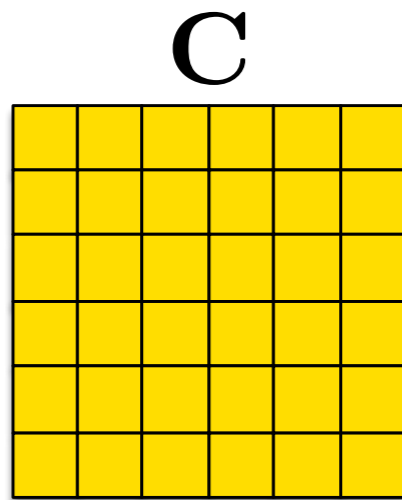
Matrix Multiplication Expressed Functionally



$A \gg \text{map}(\lambda \text{rowOfA} \mapsto$
 $B \gg \text{map}(\lambda \text{colOfB} \mapsto$

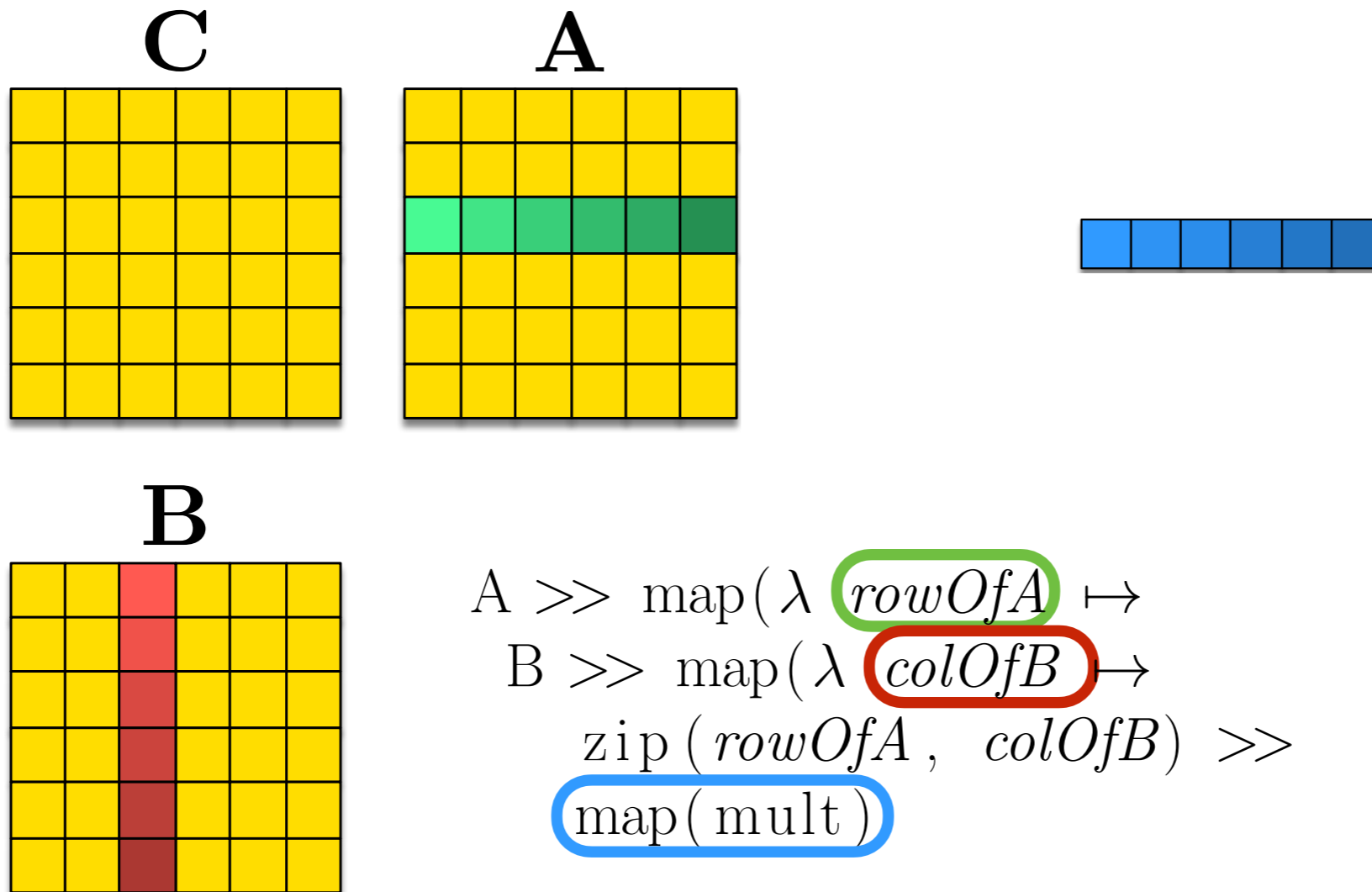
))

Matrix Multiplication Expressed Functionally



```
A >> map(λ rowOfA ↦
  B >> map(λ colOfB ↦
    zip(rowOfA, colOfB) >>
      ) )
```

Matrix Multiplication Expressed Functionally

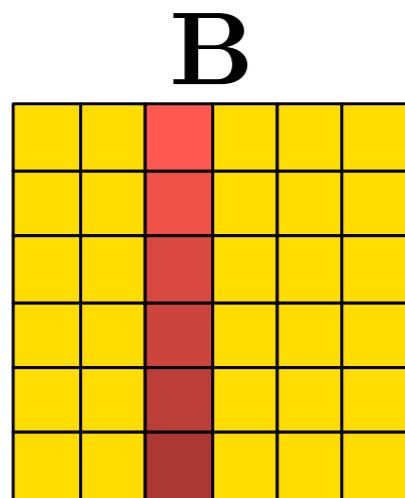
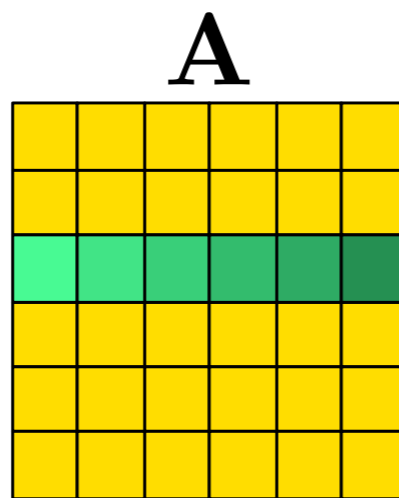
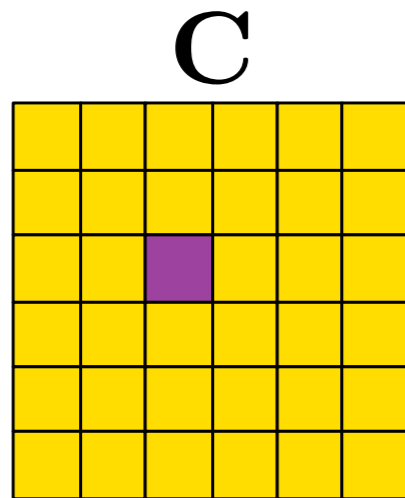


```

A >> map(λ rowOfA ↦
B >> map(λ colOfB ↦
zip(rowOfA, colOfB) >>
map(mult) ) )

```

Matrix Multiplication Expressed Functionally



```
A >> map(λ rowOfA ↦  
B >> map(λ colOfB ↦  
zip(rowOfA, colOfB) >>  
map(mult) >> reduce(0.0f, add) ) )
```

```
A >> map(λ rowOfA ↦  
  B >> map(λ colOfB ↦  
    zip(rowOfA, colOfB) >>  
    map(mult) >> reduce(0.0f, add)  
  )  
)
```

```

A >> map(λ rowOfA ↦
  B >> map(λ colOfB ↦
    zip(rowOfA, colOfB) >>
    map(mult) >> reduce(0.0f, add)
  )
)

```

Naive code
generation



```

1 for (int i = 0; i < M; i++) {
2   for (int j = 0; j < N; j++) {
3     for (int k = 0; k < K; k++) {
4       temp[k + K*N*i + K*j] =
5         mult(A[k + K*i], B[k + K*j]);
6     }
7     for (int k = 0; k < K; k++) {
8       C[j + N*i] +=
9         temp[k + K*N*i + K*j];
10    }
11  }
12 }

```

```

A >> map( $\lambda$  rowOfA  $\mapsto$ 
  B >> map( $\lambda$  colOfB  $\mapsto$ 
    zip(rowOfA, colOfB) >>
    map(mult) >> reduce(0.0f, add)
  )
)

```

Naive code
generation



```

1 for (int i = 0; i < M; i++) {
2   for (int j = 0; j < N; j++) {
3     for (int k = 0; k < K; k++) {
4       temp[k + K*N*i + K*j] =
5         mult(A[k + K*i], B[k + K*j]);
6     }
7     for (int k = 0; k < K; k++) {
8       C[j + N*i] +=
9         temp[k + K*N*i + K*j];
10    }
11  }
12 }

```



```

A >> map(λ rowOfA ↦
  B >> map(λ colOfB ↦
    zip(rowOfA, colOfB) >>
    map(mult) >> reduce(0.0f, add)
  )
)

```

Naive code
generation



```

1 for (int i = 0; i < M; i++) {
2   for (int j = 0; j < N; j++) {
3     for (int k = 0; k < K; k++) {
4       temp[k + K*N*i + K*j] =
5         mult(A[k + K*i], B[k + K*j]);
6     }
7     for (int k = 0; k < K; k++) {
8       C[j + N*i] +=
9         temp[k + K*N*i + K*j];
10    }
11  }
12 }

```

```

A >> map(λ rowOfA ↦
  B >> map(λ colOfB ↦
    zip(rowOfA, colOfB) >>
      map(mult) >> reduce(0.0f, add)
  )
)

```

Naive code
generation



```

1 for (int i = 0; i < M; i++) {
2   for (int j = 0; j < N; j++) {
3     for (int k = 0; k < K; k++) {
4       temp[k + K*N*i + K*j] =
5         mult(A[k + K*i], B[k + K*j]);
6     }
7     for (int k = 0; k < K; k++) {
8       C[j + N*i] +=
9         temp[k + K*N*i + K*j];
10    }
11  }
12 }

```

```

A >> map(λ rowOfA ↦
  B >> map(λ colOfB ↦
    zip(rowOfA, colOfB) >>
      map(mult) >> reduce(0.0f, add)
  )
)

```

Naive code
generation

How to generate
high performance code?

```

1 for (int i = 0; i < M; i++) {
2   for (int j = 0; j < N; j++) {
3     for (int k = 0; k < K; k++) {
4       temp[k + K*N*i + K*j] =
5         mult(A[k + K*i], B[k + K*j]);
6     }
7     for (int k = 0; k < K; k++) {
8       C[j + N*i] +=
9         temp[k + K*N*i + K*j];
10    }
11  }
12 }

```

?

Algorithmic Rewrite Rules

Split-join rule:

$$\text{map}(f) \implies \text{split}(k) \gg \text{map}(\text{map}(f)) \gg \text{join}$$

Map fusion rule:

$$\text{map}(f) \gg \text{map}(g) \implies \text{map}(f \gg g)$$

Map-reduce fusion rule:

$$\begin{aligned} \text{mapSeq}(f) \gg \text{reduceSeq}(z, \oplus) &\implies \\ \text{reduceSeq}(z, \lambda (acc, x) \mapsto \oplus(acc, f(x))) & \end{aligned}$$

Algorithmic Rewrite Rules

Split-join rule:

$$\text{map}(f) \implies \text{split}(k) \gg \text{map}(\text{map}(f)) \gg \text{join}$$

Map fusion rule:

$$\text{map}(f) \gg \text{map}(g) \implies \text{map}(f \gg g)$$

Map-reduce fusion rule:

$$\begin{aligned} \text{mapSeq}(f) \gg \text{reduceSeq}(z, \oplus) &\implies \\ \text{reduceSeq}(z, \lambda (acc, x) \mapsto \oplus(acc, f(x))) & \end{aligned}$$

- *Rewrite rules* are provably correct
- Express algorithmic and optimisation choices

How to apply
rewrite rules?

```

A >> map( $\lambda$  rowOfA  $\mapsto$ 
  B >> map( $\lambda$  colOfB  $\mapsto$ 
    zip(rowOfA, colOfB) >>
    map(mult) >> reduce(0.0f, add)
  )
)

```



```

1 for (int i = 0; i < M; i++) {
2   for (int j = 0; j < N; j++) {
3     for (int k = 0; k < K; k++) {
4       temp[k + K*N*i + K*j] =
5         mult(A[k + K*i], B[k + K*j]);
6     }
7   for (int k = 0; k < K; k++) {
8     C[j + N*i] +=
9       temp[k + K*N*i + K*j];
10  }
11 }
12 }

```

```

A >> map(λ rowOfA ↦
  B >> map(λ colOfB ↦
    zip(rowOfA, colOfB) >>
    map(mult) >> reduce(0.0 f, add)
  )
)

```

```

1 for (int i = 0; i < M; i++) {
2   for (int j = 0; j < N; j++) {
3     for (int k = 0; k < K; k++) {
4       temp[k + K*N i + K*j] =
5         mult(A[k + K*i], B[k + K*j]);
6     }
7   for (int k = 0; k < K; k++) {
8     C[j + N*i] +=
9       temp[k + K*N*i + K*j];
10  }
11 }
12 }

```

`map(f) ⇒ split(k) >> map(map(f)) >> join`


```

A >> map(λ rowOfA ↦
  B >> map(λ colOfB ↦
    zip(rowOfA, colOfB) >>
    map(mult) >> reduce(0.0 f, add)
  )
)

```

```

1 for (int i = 0; i < M; i++) {
2   for (int j = 0; j < N; j++) {
3     for (int k = 0; k < K; k++) {
4       temp[k + K*N i + K*j] =
5         mult(A[k + K*i], B[k + K*j]);
6     }
7   }
8   for (int k = 0; k < K; k++) {
9     C[j + N*i] +=
10      temp[k + K*N*i + K*j];
11   }
12 }

```

map(f) ⇒ split(k) >> map(map(f)) >> join

```

A >> split(m) >> map(λ rowsOfA ↦
  rowsOfA >> map(λ rowOfA ↦
    B >> map(λ colOfB ↦
      zip(rowOfA, colOfB) >>
      map(mult) >> reduce(0.0 f, add)
    )
  )
) >> join

```

```

A >> map(λ rowOfA ↦
  B >> map(λ colOfB ↦
    zip(rowOfA, colOfB) >>
    map(mult) >> reduce(0.0 f, add)
  )
)

```

```

1 for (int i = 0; i < M; i++) {
2   for (int j = 0; j < N; j++) {
3     for (int k = 0; k < K; k++) {
4       temp[k + K*N*i + K*j] =
5         mult(A[k + K*i], B[k + K*j]);
6     }
7   }
8   for (int k = 0; k < K; k++) {
9     C[j + N*i] +=
10      temp[k + K*N*i + K*j];
11   }
12 }

```

map(f) ⇒ split(k) >> map(map(f)) >> join

```

A >> split(m) >> map(λ rowsOfA ↦
  rowsOfA >> map(λ rowOfA ↦
    B >> map(λ colOfB ↦
      zip(rowOfA, colOfB) >>
      map(mult) >> reduce(0.0 f, add)
    )
  )
) >> join

```

```

1 for (int i = 0; i < M/2; i++) {
2   for (int l = 0; l < 2; l++) {
3     for (int j = 0; j < N; j++) {
4       for (int k = 0; k < K; k++) {
5         temp[k + 2*K*N*i + K*N*l + K*j] =
6           mult(A[k + K*l + 2*K*i], B[k + K*j]);
7       }
8     }
9     for (int k = 0; k < K; k++) {
10      C[j + N*l + 2*N*i] +=
11       temp[k + 2*K*N*i + K*N*l + K*j];
12    }
13  }
14 }

```

```

A >> split(m) >> map(λ rowsOfA ↦
  rowsOfA >> map(λ rowOfA ↦
    B >> map(λ colOfB ↦
      zip(rowOfA, colOfB) >>
      map(mult) >> reduce(0.0f, add)
    )
  )
) >> join

```



```

1 for (int i = 0; i < M/2; i++) {
2   for (int l = 0; l < 2; l++) {
3     for (int j = 0; j < N; j++) {
4       for (int k = 0; k < K; k++) {
5         temp[k + 2*K*N*i + K*N*l + K*j] =
6           mult(A[k + K*l + 2*K*i], B[k + K*j]);
7       }
8     }
9     for (int k = 0; k < K; k++) {
10      C[j + N*l + 2*N*i] +=
11        temp[k + 2*K*N*i + K*N*l + K*j];
12    }
13  }
14 }

```

```

A >> split(m) >> map( $\lambda$  rowsOfA  $\mapsto$ 
rowsOfA >> map( $\lambda$  rowOfA  $\mapsto$ 
  B >> map( $\lambda$  colOfB  $\mapsto$ 
    zip(rowOfA, colOfB) >>
    map(mult) >> reduce(0.0f, add)
  )
)
) >> join

```



```

1 for (int i = 0; i < M/2; i++) {
2   for (int l = 0; l < 2; l++) {
3     for (int j = 0; j < N; j++) {
4       for (int k = 0; k < K; k++) {
5         temp[k + 2*K*N*i + K*N*l + K*j] =
6           mult(A[k + K*l + 2*K*i], B[k + K*j]);
7       }
8     }
9     for (int k = 0; k < K; k++) {
10      C[j + N*l + 2*N*i] +=
11        temp[k + 2*K*N*i + K*N*l + K*j];
12    }
13  }
14 }

```

$$X \gg \text{map}(\lambda x \mapsto Y \gg \text{map}(\lambda y \mapsto f)) \implies Y \gg \text{map}(\lambda y \mapsto X \gg \text{map}(\lambda x \mapsto f)) \gg \text{transpose}$$

```

A >> split(m) >> map(λ rowsOfA ↦
  rowsOfA >> map(λ rowOfA ↦
    B >> map(λ colOfB ↦
      zip(rowOfA, colOfB) >>
      map(mult) >> reduce(0.0f, add)
    )
  )
) >> join

```



```

1 for (int i = 0; i < M/2; i++) {
2   for (int l = 0; l < 2; l++) {
3     for (int j = 0; j < N; j++) {
4       for (int k = 0; k < K; k++) {
5         temp[k + 2*K*N*i + K*N*l + K*j] =
6           mult(A[k + K*l + 2*K*i], B[k + K*j]);
7       }
8     }
9     for (int k = 0; k < K; k++) {
10      C[j + N*l + 2*N*i] +=
11        temp[k + 2*K*N*i + K*N*l + K*j];
12    }
13  }
14 }

```

$X \gg \text{map}(\lambda x \mapsto Y \gg \text{map}(\lambda y \mapsto f))$
 \implies
 $Y \gg \text{map}(\lambda y \mapsto X \gg \text{map}(\lambda x \mapsto f)) \gg \text{transpose}$

```

A >> split(m) >> map(λ rowsOfA ↦
  B >> map(λ colOfB ↦
    rowsOfA >> map(λ rowOfA ↦
      zip(rowOfA, colOfB) >>
      map(mult) >> reduce(0.0f, add)
    )
  )
) >> transpose
) >> join

```

```

A >> split(m) >> map(λ rowsOfA ↦
  rowsOfA >> map(λ rowOfA ↦
    B >> map(λ colOfB ↦
      zip(rowOfA, colOfB) >>
      map(mult) >> reduce(0.0f, add)
    )
  )
) >> join

```



```

1 for (int i = 0; i < M/2; i++) {
2   for (int l = 0; l < 2; l++) {
3     for (int j = 0; j < N; j++) {
4       for (int k = 0; k < K; k++) {
5         temp[k + 2*K*N*i + K*N*l + K*j] =
6           mult(A[k + K*l + 2*K*i], B[k + K*j]);
7       }
8       for (int k = 0; k < K; k++) {
9         C[j + N*l + 2*N*i] +=
10          temp[k + 2*K*N*i + K*N*l + K*j];
11      }
12    }
13  }
14 }

```

$X \gg \text{map}(\lambda x \mapsto Y \gg \text{map}(\lambda y \mapsto f))$
 \implies
 $Y \gg \text{map}(\lambda y \mapsto X \gg \text{map}(\lambda x \mapsto f)) \gg \text{transpose}$

```

A >> split(m) >> map(λ rowsOfA ↦
  B >> map(λ colOfB ↦
    rowsOfA >> map(λ rowOfA ↦
      zip(rowOfA, colOfB) >>
      map(mult) >> reduce(0.0f, add)
    )
  )
) >> transpose
) >> join

```



```

1 for (int i = 0; i < M/2; i++) {
2   for (int j = 0; j < N; j++) {
3     for (int l = 0; l < 2; l++) {
4       for (int k = 0; k < K; k++) {
5         temp[k + 2*K*N*i + K*N*l + K*j] =
6           mult(A[k + K*l + 2*K*i], B[k + K*j]);
7       }
8       for (int k = 0; k < K; k++) {
9         C[j + N*l + 2*N*i] +=
10          temp[k + 2*K*N*i + K*N*l + K*j];
11      }
12    }
13  }
14 }

```

```

A >> split(m) >> map( $\lambda$  rowsOfA  $\mapsto$ 
  B >> map( $\lambda$  colOfB  $\mapsto$ 
    rowsOfA >> map( $\lambda$  rowOfA  $\mapsto$ 
      zip(rowOfA, colOfB) >>
      map(mult) >> reduce(0.0f, add)
    )
  ) >> transpose
) >> join

```



```

1 for (int i = 0; i < M/2; i++) {
2   for (int j = 0; j < N; j++) {
3     for (int l = 0; l < 2; l++) {
4       for (int k = 0; k < K; k++) {
5         temp[k + 2*K*N*i + K*N*l + K*j] =
6           mult(A[k + K*l + 2*K*i], B[k + K*j]);
7       }
8     }
9     for (int k = 0; k < K; k++) {
10      C[j + N*l + 2*N*i] +=
11        temp[k + 2*K*N*i + K*N*l + K*j];
12    }
13  }
14 }

```

```

A >> split(m) >> map( $\lambda$  rowsOfA  $\mapsto$ 
  B >> map( $\lambda$  colOfB  $\mapsto$ 
    rowsOfA >> map( $\lambda$  rowOfA  $\mapsto$ 
      zip(rowOfA, colOfB) >>
      map(mult) >> reduce(0.0f, add)
    )
  ) >> transpose
) >> join

```



```

1 for (int i = 0; i < M/2; i++) {
2   for (int j = 0; j < N; j++) {
3     for (int l = 0; l < 2; l++) {
4       for (int k = 0; k < K; k++) {
5         temp[k + 2*K*N*i + K*N*l + K*j] =
6           mult(A[k + K*l + 2*K*i], B[k + K*j]);
7       }
8       for (int k = 0; k < K; k++) {
9         C[j + N*l + 2*N*i] +=
10          temp[k + 2*K*N*i + K*N*l + K*j];
11      }
12    }
13  }
14 }

```

$\text{map}(f) \implies \text{split}(k) \gg \text{map}(\text{map}(f)) \gg \text{join}$


```

A >> split(m) >> map( $\lambda$  rowsOfA  $\mapsto$ 
  B >> map( $\lambda$  colOfB  $\mapsto$ 
    rowsOfA >> map( $\lambda$  rowOfA  $\mapsto$ 
      zip(rowOfA, colOfB) >>
      map(mult) >> reduce(0.0f, add)
    )
  ) >> transpose
) >> join

```

```

1 for (int i = 0; i < M/2; i++) {
2   for (int j = 0; j < N; j++) {
3     for (int l = 0; l < 2; l++) {
4       for (int k = 0; k < K; k++) {
5         temp[k + 2*K*N*i + K*N*l + K*j] =
6           mult(A[k + K*l + 2*K*i], B[k + K*j]);
7       }
8       for (int k = 0; k < K; k++) {
9         C[j + N*l + 2*N*i] +=
10          temp[k + 2*K*N*i + K*N*l + K*j];
11       }
12     }
13   }
14 }

```

$\text{map}(f) \implies \text{split}(k) \gg \text{map}(\text{map}(f)) \gg \text{join}$

```

A >> split(m) >> map( $\lambda$  rowsOfA  $\mapsto$ 
  B >> split(n) >> map( $\lambda$  colsOfB  $\mapsto$ 
    colsOfB >> map( $\lambda$  colOfB  $\mapsto$ 
      rowsOfA >> map( $\lambda$  rowOfA  $\mapsto$ 
        zip(rowOfA, colOfB) >>
        map(mult) >> reduce(0.0f, add)
      )
    )
  ) >> join >> transpose
) >> join

```

```

A >> split(m) >> map(λ rowsOfA ↦
  B >> map(λ colOfB ↦
    rowsOfA >> map(λ rowOfA ↦
      zip(rowOfA, colOfB) >>
      map(mult) >> reduce(0.0f, add)
    )
  ) >> transpose
) >> join

```

```

1 for (int i = 0; i < M/2; i++) {
2   for (int j = 0; j < N; j++) {
3     for (int l = 0; l < 2; l++) {
4       for (int k = 0; k < K; k++) {
5         temp[k + 2*K*N*i + K*N*l + K*j] =
6           mult(A[k + K*l + 2*K*i], B[k + K*j]);
7       }
8       for (int k = 0; k < K; k++) {
9         C[j + N*l + 2*N*i] +=
10          temp[k + 2*K*N*i + K*N*l + K*j];
11      }
12    }
13  }
14 }

```

map(f) ⇒ split(k) >> map(map(f)) >> join

```

A >> split(m) >> map(λ rowsOfA ↦
  B >> split(n) >> map(λ colsOfB ↦
    colsOfB >> map(λ colOfB ↦
      rowsOfA >> map(λ rowOfA ↦
        zip(rowOfA, colOfB) >>
        map(mult) >> reduce(0.0f, add)
      )
    )
  ) >> join >> transpose
) >> join

```

```

1 for (int i = 0; i < M/2; i++) {
2   for (int j = 0; j < N/2; j++) {
3     for (int m = 0; m < 2; m++) {
4       for (int l = 0; l < 2; l++) {
5         for (int k = 0; k < K; k++) {
6           temp[k + 4*K*N*i + 2*K*N*l + 2*K*j
7             + K*m] =
8             mult(A[k + K*l + 2*K*i], B[k + K*
9               m + 2*K*j]);
10          }
11          for (int k = 0; k < K; k++) {
12            C[m + 2*j + 2*N*l + 4*N*i] +=
13              temp[k + 4*K*N*i + 2*K*N*l + 2*
14                K*j + K*m];
15          }
16        }
17      }
18    }
19  }

```

```

A >> split(m) >> map(λ rowsOfA ↦
  B >> split(n) >> map(λ colsOfB ↦
    colsOfB >> map(λ colOfB ↦
      rowsOfA >> map(λ rowOfA ↦
        zip(rowOfA, colOfB) >>
          map(mult) >> reduce(0.0f, add) ) )
    ) >> join >> transpose
  ) >> join

```



```

1 for (int i = 0; i < M/2; i++) {
2   for (int j = 0; j < N/2; j++) {
3     for (int m = 0; m < 2; m++) {
4       for (int l = 0; l < 2; l++) {
5         for (int k = 0; k < K; k++) {
6           temp[k + 4*K*N*i + 2*K*N*l + 2*K*j
7             + K*m] =
8             mult(A[k + K*l + 2*K*i], B[k + K*
9               m + 2*K*j]);
10          }
11          for (int k = 0; k < K; k++) {
12            C[m + 2*j + 2*N*l + 4*N*i] +=
13              temp[k + 4*K*N*i + 2*K*N*l + 2*
14                K*j + K*m];
15          }
16        }
17      }
18    }
19  }

```

```

A >> split(m) >> map(λ rowsOfA ↦
B >> split(n) >> map(λ colsOfB ↦
  colsOfB >> map(λ colOfB ↦
    rowsOfA >> map(λ rowOfA ↦
      zip(rowOfA, colOfB) >>
        map(mult) >> reduce(0.0f, add) ) )
    ) >> join >> transpose
  ) >> join

```



```

1 for (int i = 0; i < M/2; i++) {
2   for (int j = 0; j < N/2; j++) {
3     for (int m = 0; m < 2; m++) {
4       for (int l = 0; l < 2; l++) {
5         for (int k = 0; k < K; k++) {
6           temp[k + 4*K*N*i + 2*K*N*l + 2*K*j
7             + K*m] =
8             mult(A[k + K*l + 2*K*i], B[k + K*
9               m + 2*K*j]);
10          }
11          for (int k = 0; k < K; k++) {
12            C[m + 2*j + 2*N*l + 4*N*i] +=
13              temp[k + 4*K*N*i + 2*K*N*l + 2*
14                K*j + K*m];
15          }
16        }
17      }
18    }
19  }

```

$$\begin{array}{c}
X \gg \text{map}(\lambda x \mapsto Y \gg \text{map}(\lambda y \mapsto f)) \\
\implies \\
Y \gg \text{map}(\lambda y \mapsto X \gg \text{map}(\lambda x \mapsto f)) \gg \text{transpose}
\end{array}$$

```

A >> split(m) >> map(λ rowsOfA ↦
B >> split(n) >> map(λ colsOfB ↦
  colsOfB >> map(λ colOfB ↦
    rowsOfA >> map(λ rowOfA ↦
      zip(rowOfA, colOfB) >>
        map(mult) >> reduce(0.0f, add) ) )
    ) >> join >> transpose
  ) >> join

```



```

1 for (int i = 0; i < M/2; i++) {
2   for (int j = 0; j < N/2; j++) {
3     for (int m = 0; m < 2; m++) {
4       for (int l = 0; l < 2; l++) {
5         for (int k = 0; k < K; k++) {
6           temp[k + 4*K*N*i + 2*K*N*l + 2*K*j
7             + K*m] =
8             mult(A[k + K*l + 2*K*i], B[k + K*
9               m + 2*K*j]);
10          }
11          for (int k = 0; k < K; k++) {
12            C[m + 2*j + 2*N*l + 4*N*i] +=
13              temp[k + 4*K*N*i + 2*K*N*l + 2*
14                K*j + K*m];
15          }
16        }
17      }
18    }
19  }

```

$X \gg \text{map}(\lambda x \mapsto Y \gg \text{map}(\lambda y \mapsto f))$

\Rightarrow

$Y \gg \text{map}(\lambda y \mapsto X \gg \text{map}(\lambda x \mapsto f)) \gg \text{transpose}$

```

A >> split(m) >> map(λ rowsOfA ↦
B >> split(n) >> map(λ colsOfB ↦
  rowsOfA >> map(λ rowOfA ↦
    colsOfB >> map(λ colOfB ↦
      zip(rowOfA, colOfB) >>
        map(mult) >> reduce(0.0f, add) )
    ) >> transpose
  ) >> join >> transpose
  ) >> join

```

```

A >> split(m) >> map(λ rowsOfA ↦
B >> split(n) >> map(λ colsOfB ↦
  colsOfB >> map(λ colOfB ↦
    rowsOfA >> map(λ rowOfA ↦
      zip(rowOfA, colOfB) >>
        map(mult) >> reduce(0.0f, add) ) )
  ) >> join >> transpose
) >> join

```



```

1 for (int i = 0; i < M/2; i++) {
2   for (int j = 0; j < N/2; j++) {
3     for (int m = 0; m < 2; m++) {
4       for (int l = 0; l < 2; l++) {
5         for (int k = 0; k < K; k++) {
6           temp[k + 4*K*N*i + 2*K*N*l + 2*K*j
7             + K*m] =
8             mult(A[k + K*l + 2*K*i], B[k + K*
9               m + 2*K*j]);
10          }
11          for (int k = 0; k < K; k++) {
12            C[m + 2*j + 2*N*l + 4*N*i] +=
13              temp[k + 4*K*N*i + 2*K*N*l + 2*
14                K*j + K*m];
15          }
16        }
17      }
18    }
19  }

```

$X \gg \text{map}(\lambda x \mapsto Y \gg \text{map}(\lambda y \mapsto f))$

\Rightarrow

$Y \gg \text{map}(\lambda y \mapsto X \gg \text{map}(\lambda x \mapsto f)) \gg \text{transpose}$



```

A >> split(m) >> map(λ rowsOfA ↦
B >> split(n) >> map(λ colsOfB ↦
  rowsOfA >> map(λ rowOfA ↦
    colsOfB >> map(λ colOfB ↦
      zip(rowOfA, colOfB) >>
        map(mult) >> reduce(0.0f, add) )
  ) >> transpose
) >> join >> transpose
) >> join

```



```

1 for (int i = 0; i < M/2; i++) {
2   for (int j = 0; j < N/2; j++) {
3     for (int l = 0; l < 2; l++) {
4       for (int m = 0; m < 2; m++) {
5         for (int k = 0; k < K; k++) {
6           temp[k + 4*K*N*i + 2*K*N*l + 2*K*j
7             + K*m] =
8             mult(A[k + K*l + 2*K*i], B[k + K*
9               m + 2*K*j]);
10          }
11          for (int k = 0; k < K; k++) {
12            C[m + 2*j + 2*N*l + 4*N*i] +=
13              temp[k + 4*K*N*i + 2*K*N*l + 2*
14                K*j + K*m];
15          }
16        }
17      }
18    }
19  }

```

After algorithmic rewrites...

Tiled Matrix Multiplication

```
 $\lambda$  (A, B)  $\mapsto$ 
  A >> split(m) >> map( $\lambda$  nRowsOfA  $\mapsto$ 
    B >> split(n) >> map( $\lambda$  mColsOfB  $\mapsto$ 
      zip( transpose(nRowsOfA) >> split(k),
          transpose(mColsOfB) >> split(k) ) >>
      reduceSeq( init = make2DArray(n,m, 0.0 f),
         $\lambda$  (accTile, (tileOfA, tileOfB))  $\mapsto$ 
          zip(accTile, transpose(tileOfA)) >>
          map( $\lambda$  (accRow, rowOfTileOfA)  $\mapsto$ 
            zip(accRow, transpose(tileOfB)) >>
            map( $\lambda$  (acc, colOfTileOfB)  $\mapsto$ 
              zip(rowOfTileOfA, colOfTileOfB) >>
              map(mult) >> reduce(acc, add)
            ) >> join
          )
        ) >> transpose() >>
        map(transpose) >> transpose
      ) >> join >> transpose
    ) >> join
```

```
1 for (int i = 0; i < M/2; i++) {
2   for (int j = 0; j < N/2; j++) {
3     for (int k = 0; k < K/4; k++) {
4       for (int l = 0; l < 2; l++) {
5         for (int m = 0; m < 2; m++) {
6           for (int n = 0; n < 4; n++) {
7             temp[n + 4*m + 8*N*i + 16*j + 8*l] =
8               mult(
9                 A[n + 2*K*i + 4*k + K*l],
10                B[n + 2*K*j + 4*k + K*m]
11              );
12           }
13         for (int n = 0; n < 4; n++) {
14           C[m + 2*N*i + 2*j + N*l] +=
15             temp[n + 4*m + 8*N*i + 16*j + 8*l];
16         }
17       }
18     }
19   }
20 }
21 }
```


How to map to
OpenCL?

OpenCL Specific Rewrite Rules

- *Rewrite rules* express mapping and optimisation choices
- *Patterns* correspond to OpenCL concepts

Examples:

OpenCL thread hierarchy:

$$\text{map}(f) \implies \text{mapGlb}_{\{0,1,2\}}(f)$$

$$\text{map}(f) \implies \text{mapLcl}_{\{0,1,2\}}(f)$$

OpenCL memory hierarchy:

$$f \implies \text{toPrivate}(f)$$

$$f \implies \text{toLocal}(f)$$

$$f \implies \text{toGlobal}(f)$$

OpenCL vector types and operations:

$$\begin{aligned} \text{map}(f) &\implies \text{asVector}(n, b) \\ &\gg \text{map}(\text{vectorize}(n, f)) \gg \text{asScalar} \end{aligned}$$

```

λ (A, B) ↦
  A >> split(m) >> map(λ nRowsOfA ↦
    B >> split(n) >> map(λ mColsOfB ↦
      zip( transpose(nRowsOfA) >> split(k),
           transpose(mColsOfB) >> split(k) ) >>
      reduceSeq( init = make2DArray(n,m, 0.0 f),
                 λ (accTile, (tileOfA, tileOfB)) ↦
                   zip(accTile, transpose(tileOfA)) >>
                   map(λ (accRow, rowOfTileOfA) ↦
                     zip(accRow, transpose(tileOfB)) >>
                     map(λ (acc, colOfTileOfB) ↦
                       zip(rowOfTileOfA, colOfTileOfB) >>
                       map(mult) >> reduce(acc, add)
                     ) >> join
                   )
                 ) >> transpose() >>
      map(transpose) >> transpose
    ) >> join >> transpose
  ) >> join

```



```

1 for (int i = 0; i < M/2; i++) {
2   for (int j = 0; j < N/2; j++) {
3     for (int k = 0; k < K/4; k++) {
4       for (int l = 0; l < 2; l++) {
5         for (int m = 0; m < 2; m++) {
6           for (int n = 0; n < 4; n++) {
7             temp[n + 4*m + 8*N*i + 16*j + 8*l] =
8               mult(
9                 A[n + 2*K*i + 4*k + K*l],
10                B[n + 2*K*j + 4*k + K*m]
11              );
12           }
13         for (int n = 0; n < 4; n++) {
14           C[m + 2*N*i + 2*j + N*l] +=
15             temp[n + 4*m + 8*N*i + 16*j + 8*l];
16         }
17       }
18     }
19   }
20 }
21 }

```

```

λ (A, B) ↦
A >> split(m) >> map(λ nRowsOfA ↦
B >> split(n) >> map(λ mColsOfB ↦
zip( transpose(nRowsOfA) >> split(k),
      transpose(mColsOfB) >> split(k) ) >>
reduceSeq( init = make2DArray(n,m, 0.0 f),
λ (accTile, (tileOfA, tileOfB)) ↦
zip(accTile, transpose(tileOfA)) >>
map(λ (accRow, rowOfTileOfA) ↦
zip(accRow, transpose(tileOfB)) >>
map(λ (acc, colOfTileOfB) ↦
zip(rowOfTileOfA, colOfTileOfB) >>
map(mult) >> reduce(acc, add)
) >> join
)
) >> transpose() >>
map(transpose) >> transpose
) >> join >> transpose
) >> join

```



```

1 for (int i = 0; i < M/2; i++) {
2   for (int j = 0; j < N/2; j++) {
3     for (int k = 0; k < K/4; k++) {
4       for (int l = 0; l < 2; l++) {
5         for (int m = 0; m < 2; m++) {
6           for (int n = 0; n < 4; n++) {
7             temp[n + 4*m + 8*N*i + 16*j + 8*l] =
8               mult(
9                 A[n + 2*K*i + 4*k + K*l],
10                B[n + 2*K*j + 4*k + K*m]
11              );
12           }
13         for (int n = 0; n < 4; n++) {
14           C[m + 2*N*i + 2*j + N*l] +=
15             temp[n + 4*m + 8*N*i + 16*j + 8*l];
16         }
17       }
18     }
19   }
20 }
21 }

```

```

zip(a, b) >> map(f)
⇒
zip(asVector(n, a), asVector(n, b))
>> map(vectorize(n, f)) >> asScalar

```

```

λ (A, B) ↦
A >> split(m) >> map(λ nRowsOfA ↦
B >> split(n) >> map(λ mColsOfB ↦
  zip( transpose(nRowsOfA) >> split(k),
        transpose(mColsOfB) >> split(k) ) >>
reduceSeq( init = make2DArray(n,m, 0.0 f),
λ (accTile, (tileOfA, tileOfB)) ↦
  zip(accTile, transpose(tileOfA)) >>
map(λ (accRow, rowOfTileOfA) ↦
  zip(accRow, transpose(tileOfB)) >>
map(λ (acc, colOfTileOfB) ↦
  zip(rowOfTileOfA >> asVector(k),
      colOfTileOfB >> asVector(k)) >>
  map(mult4) >> asScalar >>
  reduce(acc, add)
) >> join
) >> transpose() >>
map(transpose) >> transpose
) >> join >> transpose
) >> join

```

```

1 for (int i = 0; i < M/2; i++) {
2   for (int j = 0; j < N/2; j++) {
3     for (int k = 0; k < K/4; k++) {
4       for (int l = 0; l < 2; l++) {
5         for (int m = 0; m < 2; m++) {
6           float4 t = mult4(
7             vload4(A, K*i/2 + k + K*l/4),
8             vload4(B, K*j/2 + k + K*m/4)
9           );
10          vstore4(t, temp, m + 2*N*i + 4*j + 2*l);
11          for (int n = 0; n < 4; n++) {
12            C[m + 2*N*i + 2*j + N*l] +=
13              temp[n + 4*m + 8*N*i + 16*j + 8*l];
14          }
15        }
16      }
17    }
18  }
19 }

```

```

λ (A, B) ⇨
A >> split(m) >> map(λ nRowsOfA ⇨
B >> split(n) >> map(λ mColsOfB ⇨
  zip( transpose(nRowsOfA) >> split(k),
        transpose(mColsOfB) >> split(k) ) >>
reduceSeq( init = make2DArray(n,m, 0.0 f),
  λ (accTile, (tileOfA, tileOfB)) ⇨
  zip(accTile, transpose(tileOfA)) >>
  map(λ (accRow, rowOfTileOfA) ⇨
    zip(accRow, transpose(tileOfB)) >>
    map(λ (acc, colOfTileOfB) ⇨
      zip(rowOfTileOfA >> asVector(k),
          colOfTileOfB >> asVector(k)) >>
      map(mult4) >> asScalar >>
      reduce(acc, add)
    ) >> join
  )
) >> transpose() >>
map(transpose) >> transpose
) >> join >> transpose
) >> join

```

```

1 for (int i = 0; i < M/2; i++) {
2   for (int j = 0; j < N/2; j++) {
3     for (int k = 0; k < K/4; k++) {
4       for (int l = 0; l < 2; l++) {
5         for (int m = 0; m < 2; m++) {
6           float4 t = mult4(
7             vload4(A, K*i/2 + k + K*l/4),
8             vload4(B, K*j/2 + k + K*m/4)
9           );
10          vstore4(t, temp, m + 2*N*i + 4*j + 2*l);
11          for (int n = 0; n < 4; n++) {
12            C[m + 2*N*i + 2*j + N*l] +=
13              temp[n + 4*m + 8*N*i + 16*j + 8*l];
14          }
15        }
16      }
17    }
18  }
19 }

```



$\text{map}(f) \implies \text{mapSeq}(f)$
 $\text{map}(f) \implies \text{mapGlb}_{\{0,1,2\}}(f)$
 $\text{reduce}(z, \oplus) \implies \text{reduceSeq}(z, \oplus)$



```

λ (A, B) ↦
A >> split(m) >> mapGlb0(λ nRowsOfA ↦
B >> split(n) >> mapGlb1(λ mColsOfB ↦
  zip( transpose(nRowsOfA) >> split(k),
        transpose(mColsOfB) >> split(k) ) >>
reduceSeq( init = make2DArray(n,m, 0.0 f),
  λ (accTile, (tileOfA, tileOfB)) ↦
  zip(accTile, transpose(tileOfA)) >>
  mapSeq(λ (accRow, rowOfTileOfA) ↦
  zip(accRow, transpose(tileOfB)) >>
  mapSeq(λ (acc, colOfTileOfB) ↦
  zip(rowOfTileOfA >> asVector(k),
      colOfTileOfB >> asVector(k)) >>
  mapSeq(mult4) >> asScalar >>
  reduceSeq(acc, add)
  ) >> join
  )
) >> transpose() >>
map(transpose) >> transpose
) >> join >> transpose
) >> join

```



```

1 int i = get_global_id(0);
2 int j = get_global_id(1);
3 for (int k = 0;k<K/4; k++) {
4   for (int l = 0;l<2; l++) {
5     for (int m = 0;m<2; m++) {
6       float4 t = mult4(
7         vload4(A, K*i/2 + k + K*l/4),
8         vload4(B, K*j/2 + k + K*m/4)
9       );
10      vstore4(t, temp, m + 2*N*i + 4*j + 2*l);
11      for (int n = 0;n<4; n++) {
12        C[m + 2*N*i + 2*j + N*l] +=
13          temp[n + 4*m + 8*N*i + 16*j + 8*l];
14      }
15    }
16  }
17 }

```

$\lambda (A, B) \mapsto$
 $A \gg \text{split}(m) \gg \text{mapGlb}_0(\lambda n\text{RowsOfA} \mapsto$
 $B \gg \text{split}(n) \gg \text{mapGlb}_1(\lambda m\text{ColsOfB} \mapsto$
 $\text{zip}(\text{transpose}(n\text{RowsOfA}) \gg \text{split}(k),$
 $\text{transpose}(m\text{ColsOfB}) \gg \text{split}(k)) \gg$
 $\text{reduceSeq}(\text{init} = \text{make2DArray}(n, m, 0.0f),$
 $\lambda (acc\text{Tile}, (tile\text{OfA}, tile\text{OfB})) \mapsto$
 $\text{zip}(acc\text{Tile}, \text{transpose}(tile\text{OfA})) \gg$
 $\text{mapSeq}(\lambda (acc\text{Row}, row\text{OfTileOfA}) \mapsto$
 $\text{zip}(acc\text{Row}, \text{transpose}(tile\text{OfB})) \gg$
 $\text{mapSeq}(\lambda (acc, col\text{OfTileOfB}) \mapsto$
 $\text{zip}(row\text{OfTileOfA} \gg \text{asVector}(k),$
 $col\text{OfTileOfB} \gg \text{asVector}(k)) \gg$
 $\text{mapSeq}(\text{mult4}) \gg \text{asScalar} \gg$
 $\text{reduceSeq}(acc, \text{add})$
 $) \gg \text{join}$
 $)$
 $) \gg \text{transpose}() \gg$
 $\text{map}(\text{transpose}) \gg \text{transpose}$
 $) \gg \text{join} \gg \text{transpose}$
 $) \gg \text{join}$

```

1 int i = get_global_id(0);
2 int j = get_global_id(1);
3 for (int k = 0; k < K/4; k++) {
4     for (int l = 0; l < 2; l++) {
5         for (int m = 0; m < 2; m++) {
6             float4 t = mult4(
7                 vload4(A, K*i/2 + k + K*l/4),
8                 vload4(B, K*j/2 + k + K*m/4)
9             );
10            vstore4(t, temp, m + 2*N*i + 4*j + 2*l);
11            for (int n = 0; n < 4; n++) {
12                C[m + 2*N*i + 2*j + N*l] +=
13                    temp[n + 4*m + 8*N*i + 16*j + 8*l];
14            }
15        }
16    }
17 }

```

$\dots \implies \dots \gg id$
 $f \implies \text{toPrivate}(f)$



```

λ (A, B) ↦
A >> split(m) >> mapGlb0(λ nRowsOfA ↦
B >> split(n) >> mapGlb1(λ mColsOfB ↦
zip( transpose(nRowsOfA) >> split(k),
transpose(mColsOfB) >> split(k) ) >>
reduceSeq( init = make2DArray(n,m, 0.0f) >>
toPrivate( mapSeq( mapSeq( id ) ) ) ,
λ (accTile, (tileOfA, tileOfB)) ↦
zip( accTile, transpose(tileOfA) ) >>
mapSeq( λ (accRow, rowOfTileOfA) ↦
zip( accRow, transpose(tileOfB) ) >>
mapSeq( λ (acc, colOfTileOfB) ↦
zip( rowOfTileOfA >> asVector(k),
colOfTileOfB >> asVector(k) ) >>
mapSeq( mult4 ) >> asScalar >>
reduceSeq( acc, add )
) >> join
)
) >> toGlobal( mapSeq( mapSeq( mapSeq( id ) ) )
>> transpose() >>
map( transpose ) >> transpose
) >> join >> transpose
) >> join

```



```

1 int i = get_global_id(0);
2 int j = get_global_id(1);
3
4 float4 temp_0; float4 temp_1;
5 float4 temp_2; float4 temp_3;
6 float acc_0; float acc_1;
7 float acc_2; float acc_3;
8
9 for (int k = 0; k < K/4; k++) {
10
11     temp_0 = mult4(vload4(k + K*i/2,A),
12                 vload4(k + K*j/2,B));
13     acc_0 += temp_0.s0 + temp_0.s1 +
14            temp_0.s2 + temp_0.s3;
15
16     temp_1 = mult4(vload4(k + K*i/2,A),
17                 vload4(k + K + 2*K*j/4,B));
18     acc_1 += temp_1.s0 + temp_1.s1 +
19            temp_1.s2 + temp_1.s3;
20
21     temp_2 = mult4(vload4(k + K + 2*K*i/4,A),
22                 vload4(k + K*j/2,B));
23     acc_2 += temp_2.s0 + temp_2.s1 +
24            temp_2.s2 + temp_2.s3;
25
26     temp_3 = mult4(vload4(k + K + 2*K*i/4, A),
27                 vload4(k + K + 2*K*j/4, B));
28     acc_3 += temp_3.s0 + temp_3.s1 +
29            temp_3.s2 + temp_3.s3;
30 }
31 C[2*N*i + 2*j] = id(acc_0);
32 C[1 + 2*N*i + 2*j] = id(acc_1);
33 C[N + 2*N*i + 2*j] = id(acc_2);
34 C[1 + N + 2*N*i + 2*j] = id(acc_3);

```



```

λ (A, B) ↦
A >> split(m) >> mapGlb0(λ nRowsOfA ↦
B >> split(n) >> mapGlb1(λ mColsOfB ↦
zip( transpose(nRowsOfA) >> split(k),
transpose(mColsOfB) >> split(k) ) >>
reduceSeq( init = make2DArray(n,m, 0.0f) >>
toPrivate( mapSeq( mapSeq( id ) ) ) ,
λ (accTile, (tileOfA, tileOfB)) ↦
zip( accTile, transpose(tileOfA) ) >>
mapSeq( λ (accRow, rowOfTileOfA) ↦
zip( accRow, transpose(tileOfB) ) >>
mapSeq( λ (acc, colOfTileOfB) ↦
zip( rowOfTileOfA >> asVector(k),
colOfTileOfB >> asVector(k) ) >>
mapSeq( mult4 ) >> asScalar >>
reduceSeq( acc, add )
) >> join
)
) >> toGlobal( mapSeq( mapSeq( mapSeq( id ) ) )
>> transpose() >>
map( transpose ) >> transpose
) >> join >> transpose
) >> join

```



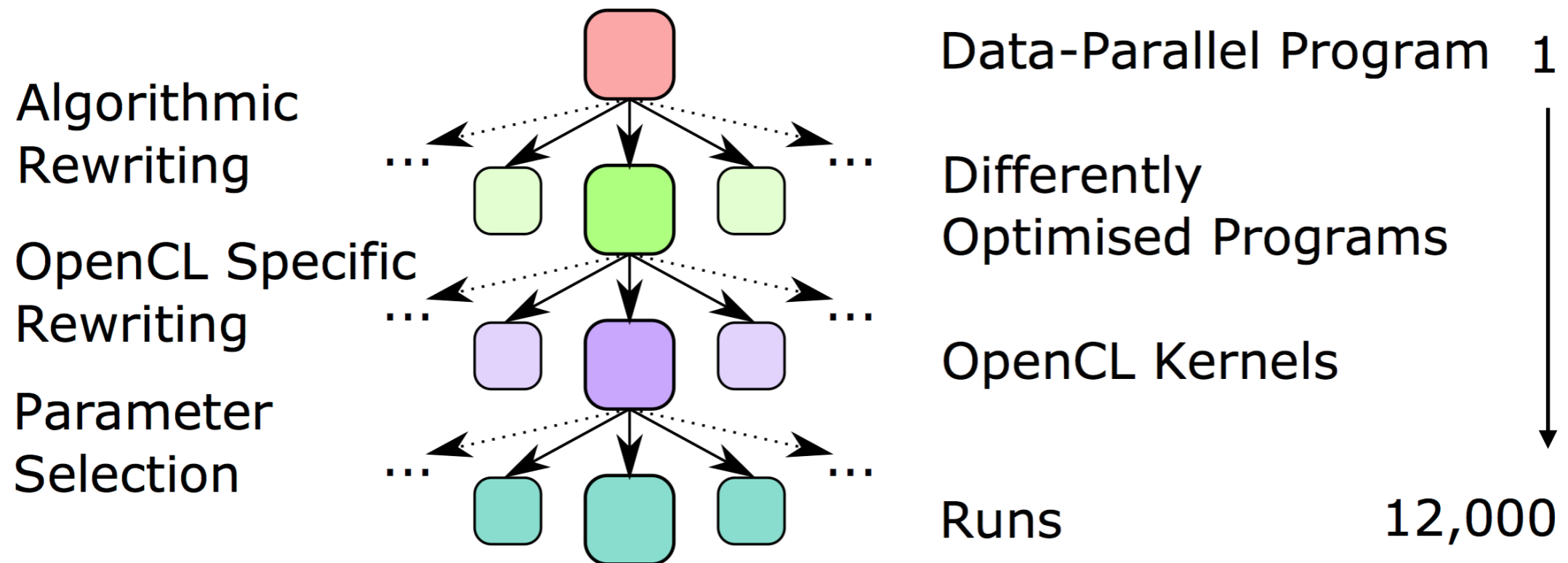
```

1 int i = get_global_id(0);
2 int j = get_global_id(1);
3
4 float4 temp_0; float4 temp_1;
5 float4 temp_2; float4 temp_3;
6 float acc_0; float acc_1;
7 float acc_2; float acc_3;
8
9 for (int k = 0; k < K/4; k++) {
10
11     temp_0 = mult4(vload4(k + K*i/2,A),
12                 vload4(k + K*j/2,B));
13     acc_0 += temp_0.s0 + temp_0.s1 +
14             temp_0.s2 + temp_0.s3;
15
16     temp_1 = mult4(vload4(k + K*i/2,A),
17                 vload4(k + K + 2*K*j/4,B));
18     acc_1 += temp_1.s0 + temp_1.s1 +
19             temp_1.s2 + temp_1.s3;
20
21     temp_2 = mult4(vload4(k + K + 2*K*i/4,A),
22                 vload4(k + K*j/2,B));
23     acc_2 += temp_2.s0 + temp_2.s1 +
24             temp_2.s2 + temp_2.s3;
25
26     temp_3 = mult4(vload4(k + K + 2*K*i/4, A),
27                 vload4(k + K + 2*K*j/4, B));
28     acc_3 += temp_3.s0 + temp_3.s1 +
29             temp_3.s2 + temp_3.s3;
30 }
31 C[2*N*i + 2*j] = id(acc_0);
32 C[1 + 2*N*i + 2*j] = id(acc_1);
33 C[N + 2*N*i + 2*j] = id(acc_2);
34 C[1 + N + 2*N*i + 2*j] = id(acc_3);

```

The generated code is highly optimised for the Mali GPU

Automated Exploration Using Rewrite Rules

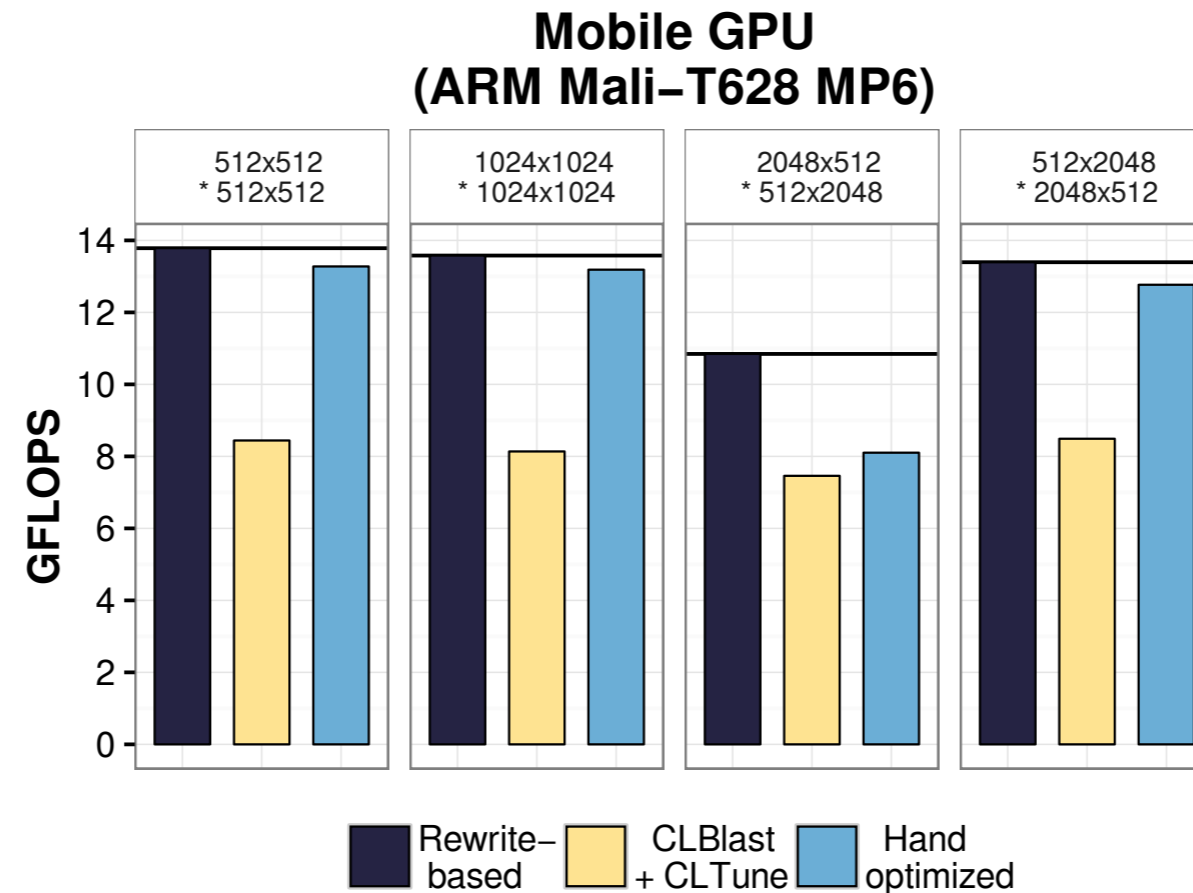


- Exploration using rewrite rules is fully automated

Automated Exploration Using Rewrite Rules

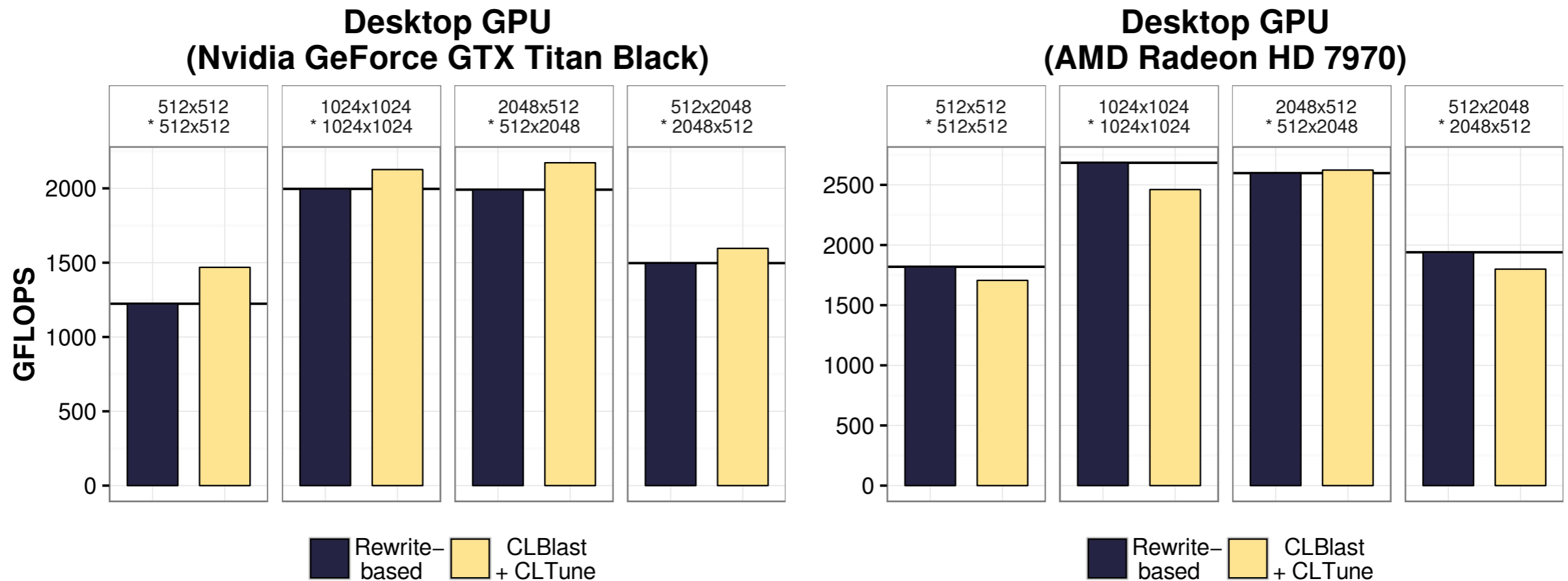
- Apply rules everywhere possible
- Stop after a certain number of applications
- Classical auto-tuning for selecting parameters
- Pick the best kernel for all sizes and devices

Mali Performance Results



- Our rewrite based approach outperforms hand optimised code on the Mali GPU

Performance Portability



- The same methodology achieves good performance across different classes of GPUs

Extensibility

- New optimisations are expressible as rewrite rules
- Automatically used in the exploration for any program

Extensibility

- New optimisations are expressible as rewrite rules
- Automatically used in the exploration for any program

Example:

OpenCL *dot* built-in rule:

```
zip(x, y) >> mapSeq(mult4) >> asScalar >> reduceSeq(z, add)
  ⇒ dot(x, y) >> reduceSeq(z, add)
```

Before:

```
1 ...
2 temp = mult4(vload4(k + K*i/2, A),
3   vload4(k + K*j/2, B));
4 acc += temp.s0 + temp.s1 +
5   temp.s2 + temp.s3;
6 ...
```

After:

```
1 ...
2 temp = dot(vload4(k + K*i/2, A),
3   vload4(k + K*j/2, B));
4 acc += temp;
5 ...
```

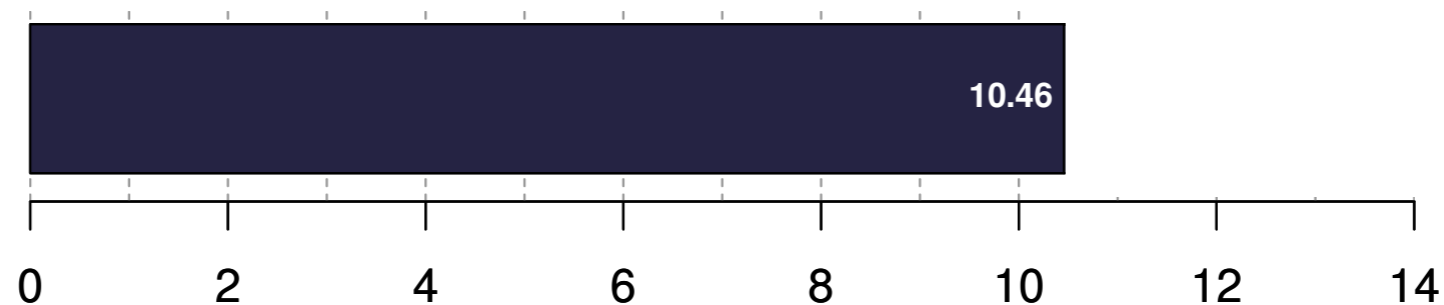

Extensibility

- New optimisations are expressible as rewrite rules
- Automatically used in the exploration for any program

Example:

OpenCL *dot* built-in rule:

```
zip(x, y) >> mapSeq(mult4) >> asScalar >> reduceSeq(z, add)  
⇒ dot(x, y) >> reduceSeq(z, add)
```



GFLOPS



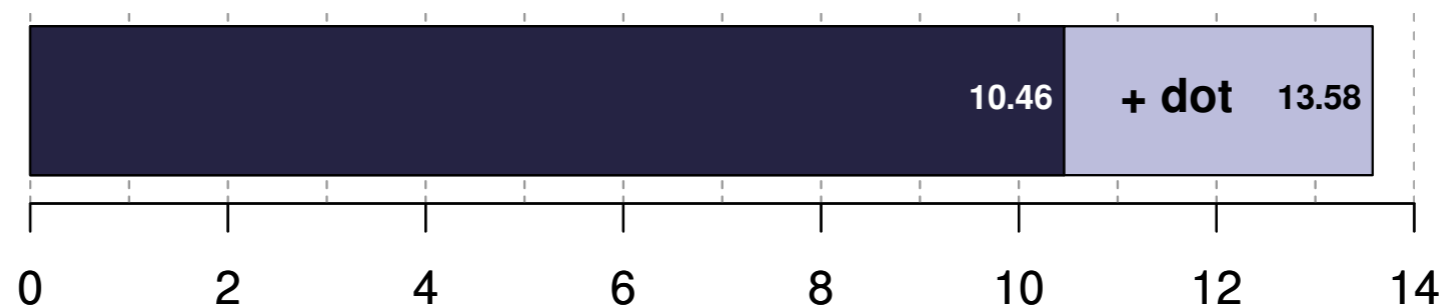
Extensibility

- New optimisations are expressible as rewrite rules
- Automatically used in the exploration for any program

Example:

OpenCL *dot* built-in rule:

```
zip(x, y) >> mapSeq(mult4) >> asScalar >> reduceSeq(z, add)
  ⇒ dot(x, y) >> reduceSeq(z, add)
```



GFLOPS



Conclusion

- Classical auto-tuning is not performance portable
- Our approach outperforms hand-tuned OpenCL code on Mali, where the auto-tuner fails to deliver
- Easily extensible by adding new rules
- Using a functional approach along with **rewrite rules** we can achieve **performance portability**

Future Work

- Building performance models to avoid executing code in the exploration to evaluate expressions
- Using reinforcement learning to guide the exploration
- Investigating a larger number of optimisations

Toomas Remmelg - toomas.remmelg@ed.ac.uk
<http://www.lift-project.org/>

Supported by:

