

# Exploiting Code Generation for Efficient LIKE Pattern Matching

**Adrian Riedl**, Philipp Fent, Maximilian Bandle, Thomas Neumann

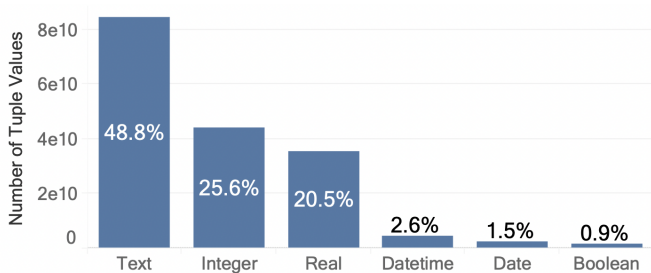
Technical University of Munich

ADMS @ VLDB 2023



# Strings are everywhere

→ Require efficient text operations



From: Vogelsgesang et al., "Get Real: How Benchmarks Fail to Represent the Real World"

- Require efficient text operations
- Should a DBMS use third-party libraries for string processing?

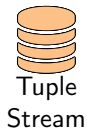
- Require efficient text operations
- Should a DBMS use third-party libraries for string processing?
- ⚠ Expensive string conversions due to DBMS specific format

- Require efficient text operations
- Should a DBMS use third-party libraries for string processing?
- ⚠ Expensive string conversions due to DBMS specific format
- ⚠ Per-tuple overhead causes significant performance loss

- Require efficient text operations
- Should a DBMS use third-party libraries for string processing?
- ⚠ Expensive string conversions due to DBMS specific format
- ⚠ Per-tuple overhead causes significant performance loss

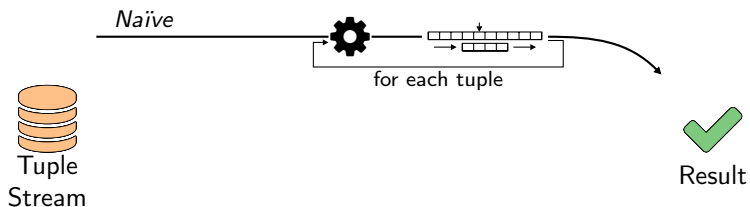
**→ We need to integrate text operations better!**

```
select count(*) from uni where name like '%TUM%';
```

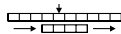


Result

```
select count(*) from uni where name like '%TUM%';
```



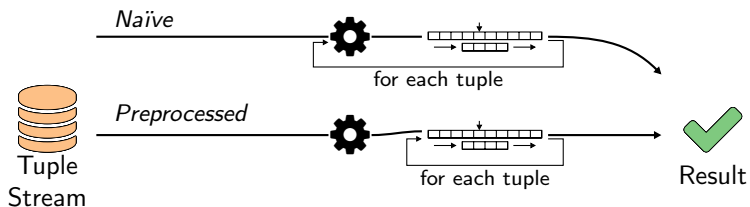
= preprocessing phase



= search phase

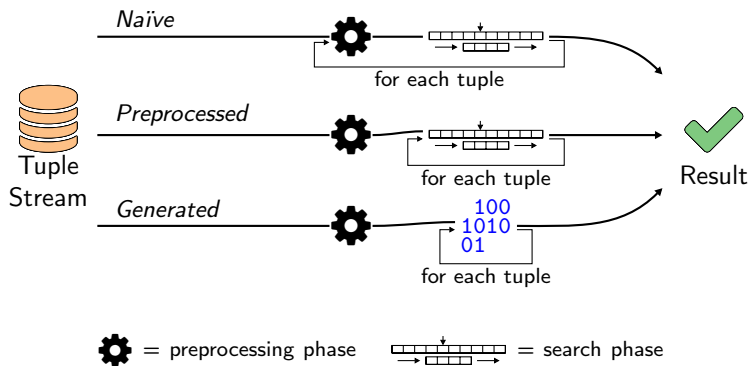


```
select count(*) from uni where name like '%TUM%';
```



 = preprocessing phase       = search phase

```
select count(*) from uni where name like '%TUM%';
```



### Naïve approach

```
KMP(text, pattern):  
    lpsTable = preprocess(pattern);  
    pPos = 0;  
    pSize = pattern.size();  
    tPos = 0;  
    tSize = text.size();  
    while (tPos - pPos + pSize <= tSize)  
        if (pattern[pPos] == text[tPos])  
            pPos++; tPos++;  
            if (pPos == pSize) return true;  
        else  
            shift = lpsTable[pPos];  
            if (shift < 0) pPos = 0; tPos++;  
            else pPos = shift;  
    return false;
```

### Naïve approach

KMP(text, pattern):

```
lpsTable = preprocess(pattern);  
pPos = 0;  
pSize = pattern.size();  
tPos = 0;  
tSize = text.size();  
while (tPos - pPos + pSize <= tSize)  
    if (pattern[pPos] == text[tPos])  
        pPos++; tPos++;  
        if (pPos == pSize) return true;  
    else  
        shift = lpsTable[pPos];  
        if (shift < 0) pPos = 0; tPos++;  
        else pPos = shift;  
return false;
```

### Naïve approach

KMP(text, pattern):

```
lpsTable = preprocess(pattern);  
pPos = 0;  
pSize = pattern.size();  
tPos = 0;  
tSize = text.size();  
while (tPos - pPos + pSize <= tSize)  
    if (pattern[pPos] == text[tPos])  
        pPos++; tPos++;  
        if (pPos == pSize) return true;  
    else  
        shift = lpsTable[pPos];  
        if (shift < 0) pPos = 0; tPos++;  
        else pPos = shift;  
return false;
```

### Naïve approach

KMP(text, pattern):

```
lpsTable = preprocess(pattern);
pPos = 0;
pSize = pattern.size();
tPos = 0;
tSize = text.size();
while (tPos - pPos + pSize <= tSize)
    if (pattern[pPos] == text[tPos])
        pPos++; tPos++;
        if (pPos == pSize) return true;
    else
        shift = lpsTable[pPos];
        if (shift < 0) pPos = 0; tPos++;
        else pPos = shift;
return false;
```

### Naïve approach

KMP(text, pattern):

```
lpsTable = preprocess(pattern);
pPos = 0;
pSize = pattern.size();
tPos = 0;
tSize = text.size();
while (tPos - pPos + pSize <= tSize)
    if (pattern[pPos] == text[tPos])
        pPos++; tPos++;
        if (pPos == pSize) return true;
    else
        shift = lpsTable[pPos];
        if (shift < 0) pPos = 0; tPos++;
        else pPos = shift;
return false;
```

### Naïve approach

KMP(text, pattern):

```
lpsTable = preprocess(pattern);
```

```
pPos = 0;
```

```
pSize = pattern.size();
```

```
tPos = 0;
```

```
tSize = text.size();
```

```
while (tPos - pPos + pSize <= tSize)
```

```
    if (pattern[pPos] == text[tPos])
```

```
        pPos++; tPos++;
```

```
        if (pPos == pSize) return true;
```

```
    else
```

```
        shift = lpsTable[pPos];
```

```
        if (shift < 0) pPos = 0; tPos++;
```

```
        else pPos = shift;
```

```
return false;
```

Extract pre-  
processing →



## From Naïve to Preprocessed

### Naïve approach

KMP(text, pattern):

```
lpsTable = preprocess(pattern);
pPos = 0;
pSize = pattern.size();
tPos = 0;
tSize = text.size();
while (tPos - pPos + pSize <= tSize)
    if (pattern[pPos] == text[tPos])
        pPos++; tPos++;
        if (pPos == pSize) return true;
    else
        shift = lpsTable[pPos];
        if (shift < 0) pPos = 0; tPos++;
        else pPos = shift;
return false;
```

Extract pre-  
processing →

### Preprocessed approach

KMP(text, pattern, lpsTable):

```
pPos = 0;
pSize = pattern.size();
tPos = 0;
tSize = text.size();
while (tPos - pPos + pSize <= tSize)
    if (pattern[pPos] == text[tPos])
        pPos++; tPos++;
        if (pPos == pSize) return true;
    else
        shift = lpsTable[pPos];
        if (shift < 0) pPos = 0; tPos++;
        else pPos = shift;
return false;
```

### Preprocessed approach

```
KMP(text, pattern, lpsTable):  
    pPos = 0;  
    pSize = pattern.size();  
    tPos = 0;  
    tSize = text.size();  
    while (tPos - pPos + pSize <= tSize)  
        if (pattern[pPos] == text[tPos])  
            pPos++; tPos++;  
            if (pPos == pSize) return true;  
        else  
            shift = lpsTable[pPos];  
            if (shift < 0) pPos = 0; tPos++;  
            else pPos = shift;  
    return false;
```

### Preprocessed approach

```
KMP(text, pattern, lpsTable):
```

```
    pPos = 0;
```

```
    pSize = pattern.size();
```

```
    tPos = 0;
```

```
    tSize = text.size();
```

```
    while (tPos - pPos + pSize <= tSize)
```

```
        if (pattern[pPos] == text[tPos]) Specialize →
```

```
            pPos++; tPos++;
```

```
            if (pPos == pSize) return true;
```

```
    else
```

```
        shift = lpsTable[pPos];
```

```
        if (shift < 0) pPos = 0; tPos++;
```

```
        else pPos = shift;
```

```
    return false;
```

### Generated approach

# Knuth-Morris-Pratt Algorithm

From Preprocessed to Generated for '%TUM%'



UMBRA



## Preprocessed approach

```
KMP(text, pattern, lpsTable):
```

```
  pPos = 0;
```

```
  pSize = pattern.size();
```

```
  tPos = 0;
```

```
  tSize = text.size();
```

```
  while (tPos - pPos + pSize <= tSize)
```

```
    if (pattern[pPos] == text[tPos])
```

```
      pPos++; tPos++;
```

```
      if (pPos == pSize) return true;
```

```
    else
```

```
      shift = lpsTable[pPos];
```

```
      if (shift < 0) pPos = 0; tPos++;
```

```
      else pPos = shift;
```

```
  return false;
```

**Specialize** →

## Generated approach

whileLoopHeader:

```
  check tPos - pPos + 3 <= text.size()
```

return false

# Knuth-Morris-Pratt Algorithm

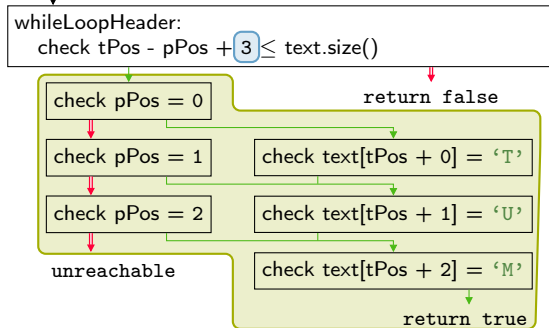
From Preprocessed to Generated for '%TUM%'

## Preprocessed approach

```
KMP(text, pattern, lpsTable):  
    pPos = 0;  
    pSize = pattern.size();  
    tPos = 0;  
    tSize = text.size();  
    while (tPos - pPos + pSize <= tSize)  
        if (pattern[pPos] == text[tPos])  
            pPos++; tPos++;  
            if (pPos == pSize) return true;  
        else  
            shift = lpsTable[pPos];  
            if (shift < 0) pPos = 0; tPos++;  
            else pPos = shift;  
    return false;
```

Specialize

## Generated approach



# Knuth-Morris-Pratt Algorithm

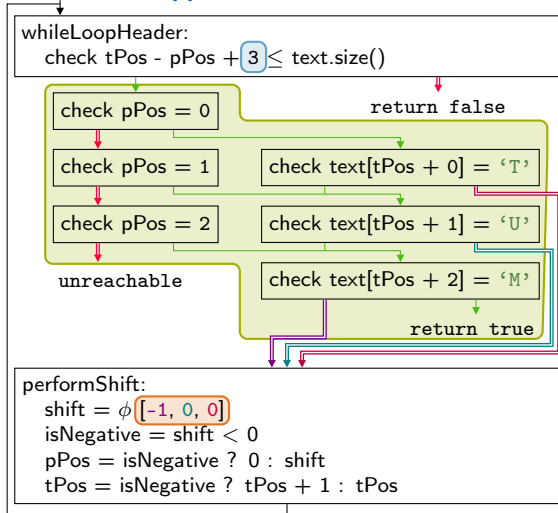
From Preprocessed to Generated for '%TUM%'

## Preprocessed approach

```
KMP(text, pattern, lpsTable):  
    pPos = 0;  
    pSize = pattern.size();  
    tPos = 0;  
    tSize = text.size();  
    while (tPos - pPos + pSize <= tSize)  
        if (pattern[pPos] == text[tPos])  
            pPos++; tPos++;  
            if (pPos == pSize) return true;  
        else  
            shift = lpsTable[pPos];  
            if (shift < 0) pPos = 0; tPos++;  
            else pPos = shift;  
    return false;
```

Specialize

## Generated approach



# Evaluation

Multi-threaded performance on Clickbench: `url like '%google%'`

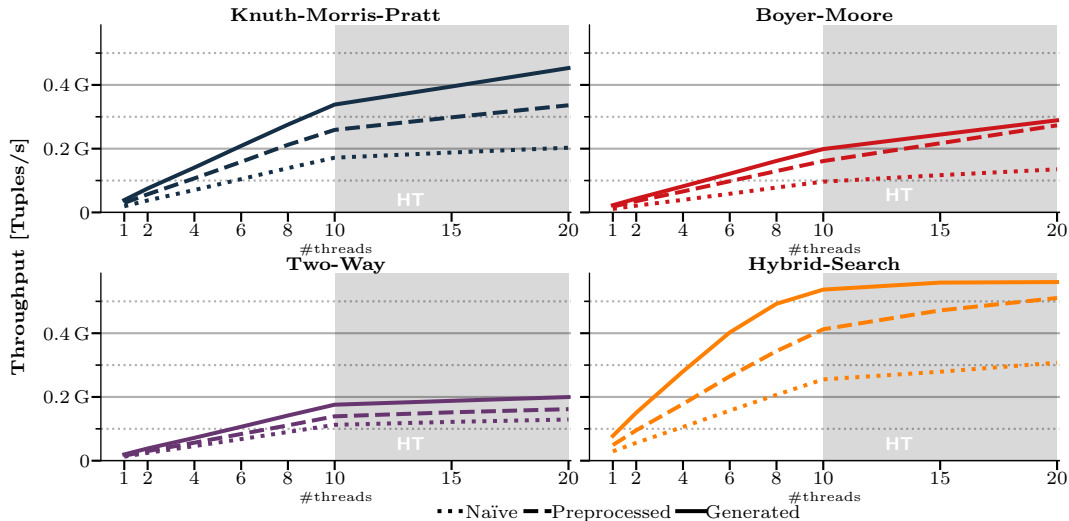
# Evaluation



UMBRA



Multi-threaded performance on Clickbench: url like '%google%'





➔ Uses SSE instruction: `pcmpistri`

```
HS(text, pattern):  
    if (pattern.size() <= 16 && text.size() >= 16)  
        iter = text.begin(), end = text.end();  
        safeMatch = 17 - pattern.size();  
        pattern16 = load16(pattern);  
        while ((iter + 16) < end)  
            match = pcmpistri(pattern16, load16(iter));  
            if (match < safeMatch) return true;  
            iter += safeMatch;  
        if (iter < end)  
            match = pcmpistri(pattern16, load16(end - 16));  
            return match < safeMatch;  
        return false;  
    return fallback(text, pattern);
```

- Uses SSE instruction: `pcmpistri`
- No preprocessing needed

```
HS(text, pattern):  
    if (pattern.size() <= 16 && text.size() >= 16)  
        iter = text.begin(), end = text.end();  
        safeMatch = 17 - pattern.size();  
        pattern16 = load16(pattern);  
        while ((iter + 16) < end)  
            match = pcmpistri(pattern16, load16(iter));  
            if (match < safeMatch) return true;  
            iter += safeMatch;  
        if (iter < end)  
            match = pcmpistri(pattern16, load16(end - 16));  
            return match < safeMatch;  
        return false;  
    return fallback(text, pattern);
```

- Uses SSE instruction: `pcmpistri`
- No preprocessing needed
- Consumes 16 bytes of input text at once

```
HS(text, pattern):  
    if (pattern.size() <= 16 && text.size() >= 16)  
        iter = text.begin(), end = text.end();  
        safeMatch = 17 - pattern.size();  
        pattern16 = load16(pattern);  
        while ((iter + 16) < end)  
            match = pcmpistri(pattern16, load16(iter));  
            if (match < safeMatch) return true;  
            iter += safeMatch;  
        if (iter < end)  
            match = pcmpistri(pattern16, load16(end - 16));  
            return match < safeMatch;  
        return false;  
    return fallback(text, pattern);
```

- Uses SSE instruction: `pcmpistri`
- No preprocessing needed
- Consumes 16 bytes of input text at once
- ✓ Generated approach is straightforward

```
HS(text, pattern):  
    if (pattern.size() <= 16 && text.size() >= 16)  
        iter = text.begin(), end = text.end();  
        safeMatch = 17 - pattern.size();  
        pattern16 = load16(pattern);  
        while ((iter + 16) < end)  
            match = pcmpistri(pattern16, load16(iter));  
            if (match < safeMatch) return true;  
            iter += safeMatch;  
        if (iter < end)  
            match = pcmpistri(pattern16, load16(end - 16));  
            return match < safeMatch;  
        return false;  
    return fallback(text, pattern);
```

- ➔ Uses SSE instruction: `pcmpistri`
- ➔ No preprocessing needed
- ➔ Consumes 16 bytes of input text at once
- ✓ Generated approach is straightforward
- ⚠ Limited to short patterns

```
HS(text, pattern):  
    if (pattern.size() <= 16 && text.size() >= 16)  
        iter = text.begin(), end = text.end();  
        safeMatch = 17 - pattern.size();  
        pattern16 = load16(pattern);  
        while ((iter + 16) < end)  
            match = pcmpistri(pattern16, load16(iter));  
            if (match < safeMatch) return true;  
            iter += safeMatch;  
        if (iter < end)  
            match = pcmpistri(pattern16, load16(end - 16));  
            return match < safeMatch;  
        return false;  
    return fallback(text, pattern);
```

- ➔ Uses SSE instruction: `pcmpistri`
- ➔ No preprocessing needed
- ➔ Consumes 16 bytes of input text at once
- ✓ Generated approach is straightforward
- ⚠ Limited to short patterns

```
HS(text, pattern):  
    if (pattern.size() <= 16 && text.size() >= 16)  
        iter = text.begin(), end = text.end();  
        safeMatch = 17 - pattern.size();  
        pattern16 = load16(pattern);  
        while ((iter + 16) < end)  
            match = pcmpistri(pattern16, load16(iter));  
            if (match < safeMatch) return true;  
            iter += safeMatch;  
        if (iter < end)  
            match = pcmpistri(pattern16, load16(end - 16));  
            return match < safeMatch;  
        return false;  
    return fallback(text, pattern);
```

➔ Code generation allows to generate code specifically for longer patterns

# SSE Search

**Example:** '%Technical University of Munich%'

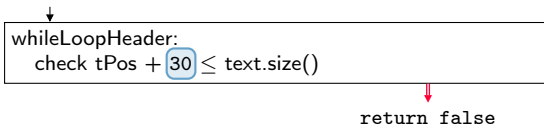
# SSE Search



UMBRA



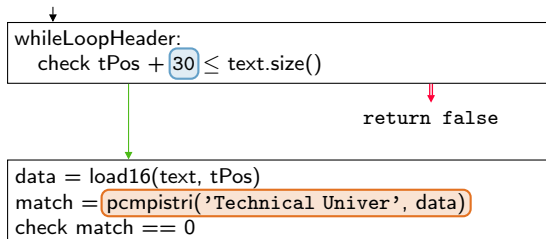
Example: '%Technical University of Munich%'





# SSE Search

Example: '%Technical University of Munich%'



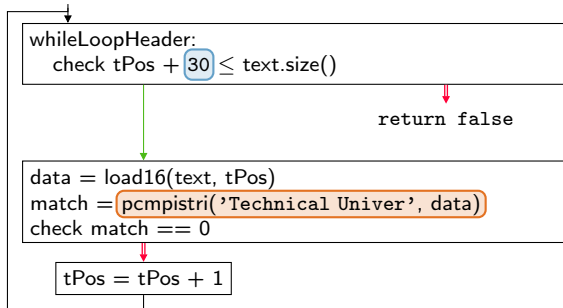
# SSE Search



UMBRA



Example: '%Technical University of Munich%'



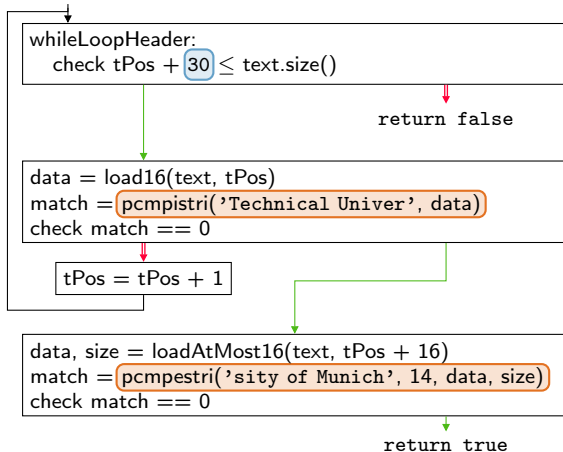
# SSE Search



UMBRA

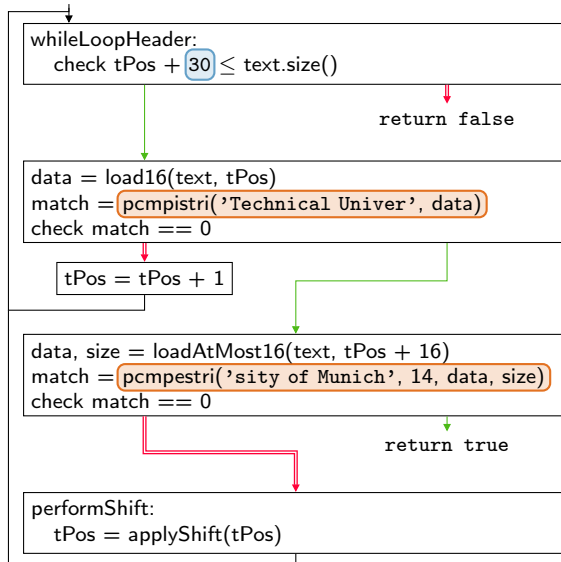


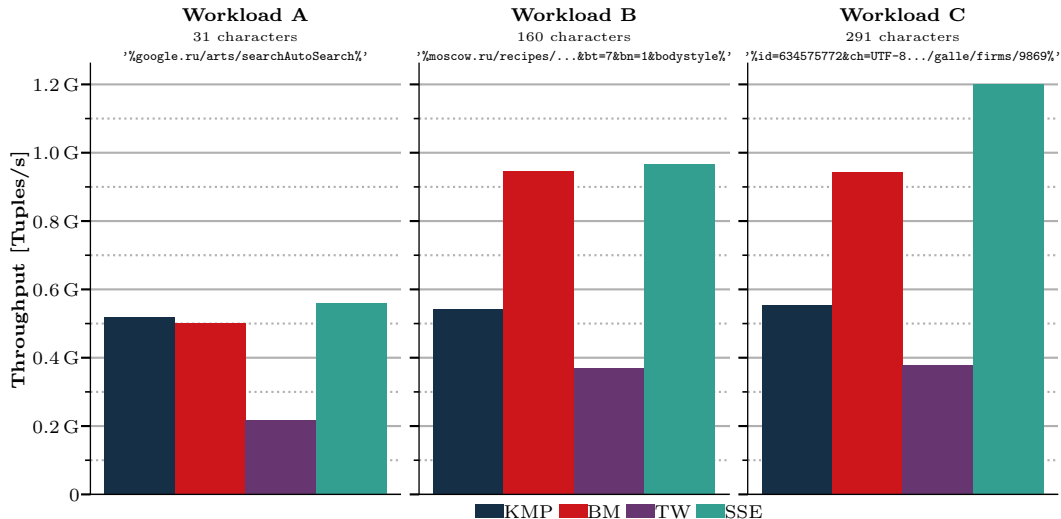
Example: '%Technical University of Munich%'



# SSE Search

Example: '%Technical University of Munich%'





# Exploiting Code Generation for Efficient LIKE Pattern Matching

- ➔ Replacing generic function calls with pattern-specific code

Adrian Riedl  
Technical University of Munich  
Adrian.Riedl@in.tum.de

## Exploiting Code Generation for Efficient LIKE Pattern Matching

- ➔ Replacing generic function calls with pattern-specific code
- ➔ Throughput improvement by up to  $2.5\times$

Adrian Riedl  
Technical University of Munich  
Adrian.Riedl@in.tum.de

## Exploiting Code Generation for Efficient LIKE Pattern Matching

- ➔ Replacing generic function calls with pattern-specific code
- ➔ Throughput improvement by up to  $2.5\times$
- ➔ Generating specialized code with SSE instructions for longer pattern

Adrian Riedl  
Technical University of Munich  
Adrian.Riedl@in.tum.de