

Masterarbeit

**Effiziente String-Verarbeitung in  
Datenbankanfragen auf hochgradig paralleler  
Hardware**

Florian Lüdiger  
Juni 2019

Gutachter:  
Prof. Dr. Jens Teubner  
Henning Funke

Technische Universität Dortmund  
Fakultät für Informatik  
Datenbanken und Informationssysteme (LS-6)  
<http://dbis.cs.tu-dortmund.de>



# Inhaltsverzeichnis

<b>1</b>	<b>Einleitung</b>	<b>1</b>
1.1	Motivation und Hintergrund . . . . .	1
1.2	Aufbau der Arbeit . . . . .	1
<b>2</b>	<b>Grundlagen des verwendeten Verfahrens</b>	<b>3</b>
2.1	Grundlagen zu Grafikkarten und deren Programmierung . . . . .	3
2.2	Das Lane-Refill-Verfahren . . . . .	3
<b>3</b>	<b>Einfacher, paralleler String-Vergleich</b>	<b>5</b>
3.1	Naive Umsetzung des Gleichheitstests . . . . .	5
3.2	Ansatzpunkte für Lane Refill . . . . .	8
3.3	Umsetzung mit Lane Refill . . . . .	8
3.4	Präfixtest als alternativer Workload . . . . .	8
3.5	Verwendete Workloads und deren Merkmale . . . . .	8
3.6	Vorstellung der Messergebnisse . . . . .	8
3.7	Diskussion der Ergebnisse . . . . .	10
<b>4</b>	<b>Paralleler Musterabgleich mit regulären Ausdrücken</b>	<b>11</b>
4.1	Grundlagen zur Verarbeitung regulärer Ausdrücke . . . . .	11
4.2	Naive Umsetzung . . . . .	11
4.3	Umsetzung mit Lane Refill . . . . .	11
4.4	Verwendete Workloads und deren Merkmale . . . . .	11
4.5	Vorstellung der Messergebnisse . . . . .	11
4.6	Diskussion der Ergebnisse . . . . .	11
<b>5</b>	<b>Ergebnis und Fazit</b>	<b>13</b>
<b>A</b>	<b>Weitere Informationen</b>	<b>15</b>
	<b>Abbildungsverzeichnis</b>	<b>17</b>
	<b>Literatur</b>	<b>19</b>



# Kapitel 1

## Einleitung

1.1 Motivation und Hintergrund

1.2 Aufbau der Arbeit



## Kapitel 2

# Grundlagen des verwendeten Verfahrens

2.1 Grundlagen zu Grafikkarten und deren Programmierung

2.2 Das Lane-Refill-Verfahren





## Kapitel 3

# Einfacher, paralleler String-Vergleich

Für die Evaluation des Lane-Refill-Verfahrens für die Verarbeitung von String-Daten wird zunächst ein einfacher String-Vergleich auf einer GPU untersucht. Ein Vergleich auf Gleichheit ist dabei die einfachste Variante von String-Verarbeitung, die vom Lane-Refill profitieren könnte. Diese Untersuchung wird dabei helfen, zu erfahren, ob die Anwendung des Lane-Refill-Verfahrens bei String-Daten allgemein Potenzial dafür bietet, den Durchsatz entsprechender Anwendungen zu erhöhen.

Zunächst wird dazu ein String-Vergleich mittels der CUDA Schnittstelle ohne spezielle Optimierungen implementiert, um einen Vergleich mit der optimierten Version durchführen zu können. Anschließend werden potenzielle Schwachstellen analysiert und daraus Ansatzpunkte für Optimierungen durch das Auffüllen leer gelaufener Lanes erarbeitet. Daraufhin wird eine Umsetzung vorgestellt, die diese Optimierungen enthält und somit durch die Verwendung geeigneter Workloads mit dem naiven Verfahren verglichen werden kann. Schließlich werden die erhaltenen Messergebnisse vorgestellt und analysiert werden.

### 3.1 Naive Umsetzung des Gleichheitstests

Als Basis für die Untersuchung wird zunächst der Gleichheitstest für Strings naiv, also ohne tiefgehende Optimierungen umgesetzt. Dies gibt Gelegenheit dazu, die Programmierung einfacher Algorithmen mithilfe der CUDA Schnittstelle für Grafikkarten darzustellen. Da die Analyse im Rahmen dieser Arbeit innerhalb einer Pipelining-Umgebung durchgeführt werden, lassen sich hier außerdem einige Besonderheiten der Implementierung erläutern.

```

1  __global__
2  void naiveKernel(
3  int *char_offset,          // indices of the first letter of every string
4  char *data_content,        // concatenated list of compare strings
5  char *search_string,       // string that will be searched for
6  int search_length,         // length of the search string
7  int line_count,           // number of lines in the data set
8  int *number_of_matches) { // return value for the number of matches
9
10     // global index of the current thread,
11     // used as the iterator in this case
12     unsigned loop_var = ((blockIdx.x * blockDim.x) + threadIdx.x);
13
14     // offset for the next element to be computed
15     unsigned step = (blockDim.x * gridDim.x);
16
17     bool active = true;
18     bool flush_pipeline = 0;
19
20     while(!flush_pipeline) {
21         // element index must not be higher than line count
22         active = loop_var < line_count;
23
24         // break computation when every line is finished and therefore inactive
25         flush_pipeline = !__ballot_sync(ALL_LANES, active);
26
27         data_length = char_offset[loop_var+1] - char_offset[loop_var] - 1;
28
29         // if the lengths of the strings don't match,
30         // the string can be discarded immediately
31         if (active && data_length != search_length)
32             active = false;
33
34         int search_id = 0;
35
36         // iterate over both strings till the end
37         // or until a non-matching character has been found
38         while(__any_sync(0xFFFFFFFF, active) && search_id < search_length) {
39
40             int data_id = search_id + char_offset[loop_var];
41
42             // when strings don't match, inactivate the lane
43             if (active && data_content[data_id] != search_string[search_id])
44                 active = false;
45
46             search_id++;
47         }

```

```
48  
49     // when comparison finishes without being inactivated,  
50     // a match has been found  
51     if (active)  
52         atomicAdd(number_of_matches, 1);  
53  
54     loop_var += step;  
55 }  
56 }
```

**Listing 3.1:** Naive Implementierung des String-Vergleichs

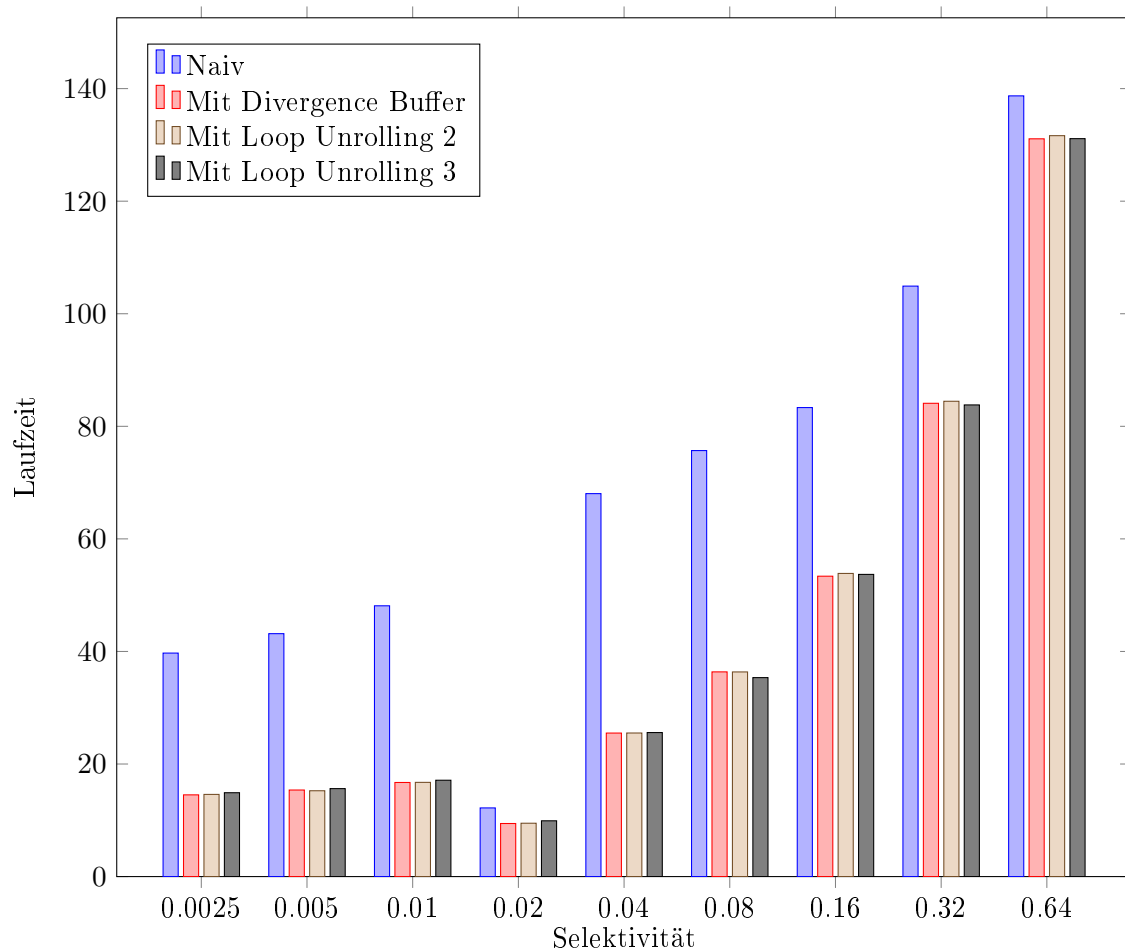
## 3.2 Ansatzpunkte für Lane Refill

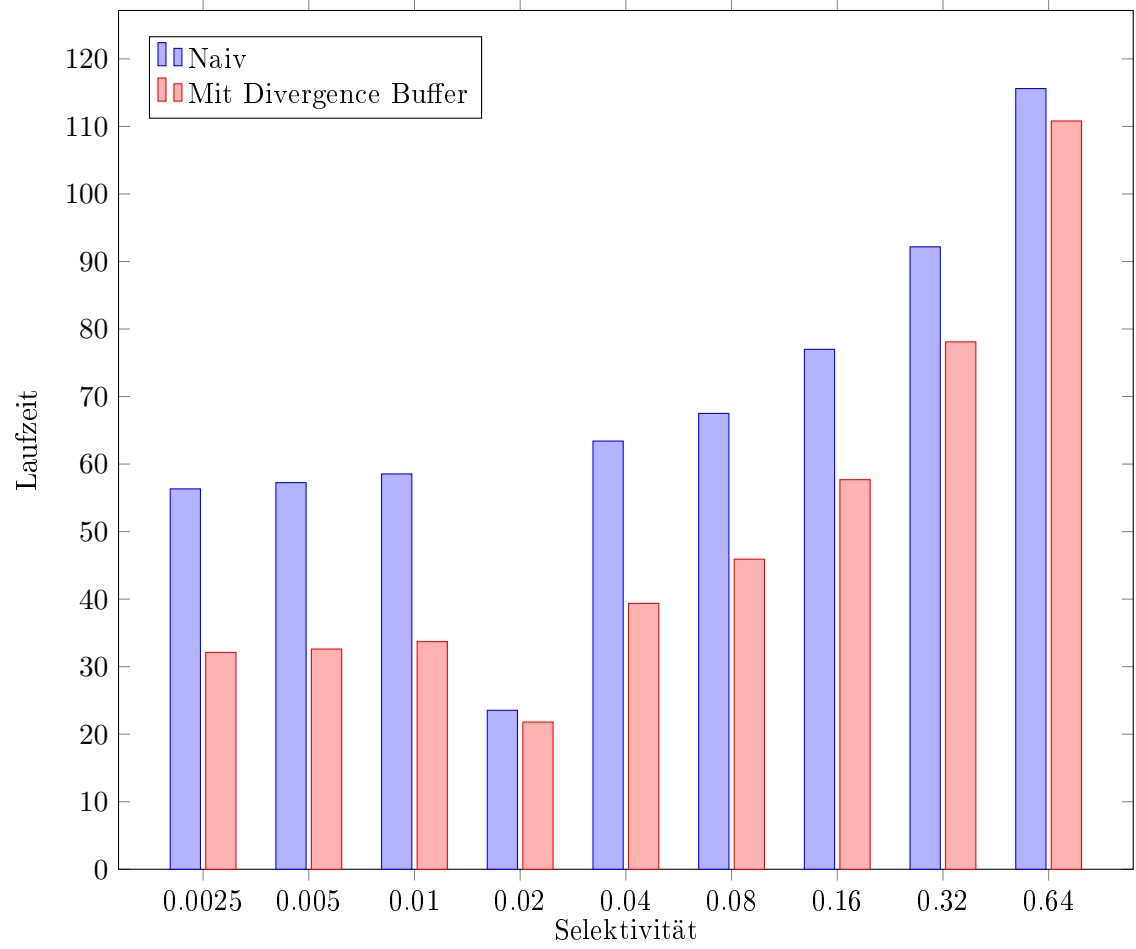
## 3.3 Umsetzung mit Lane Refill

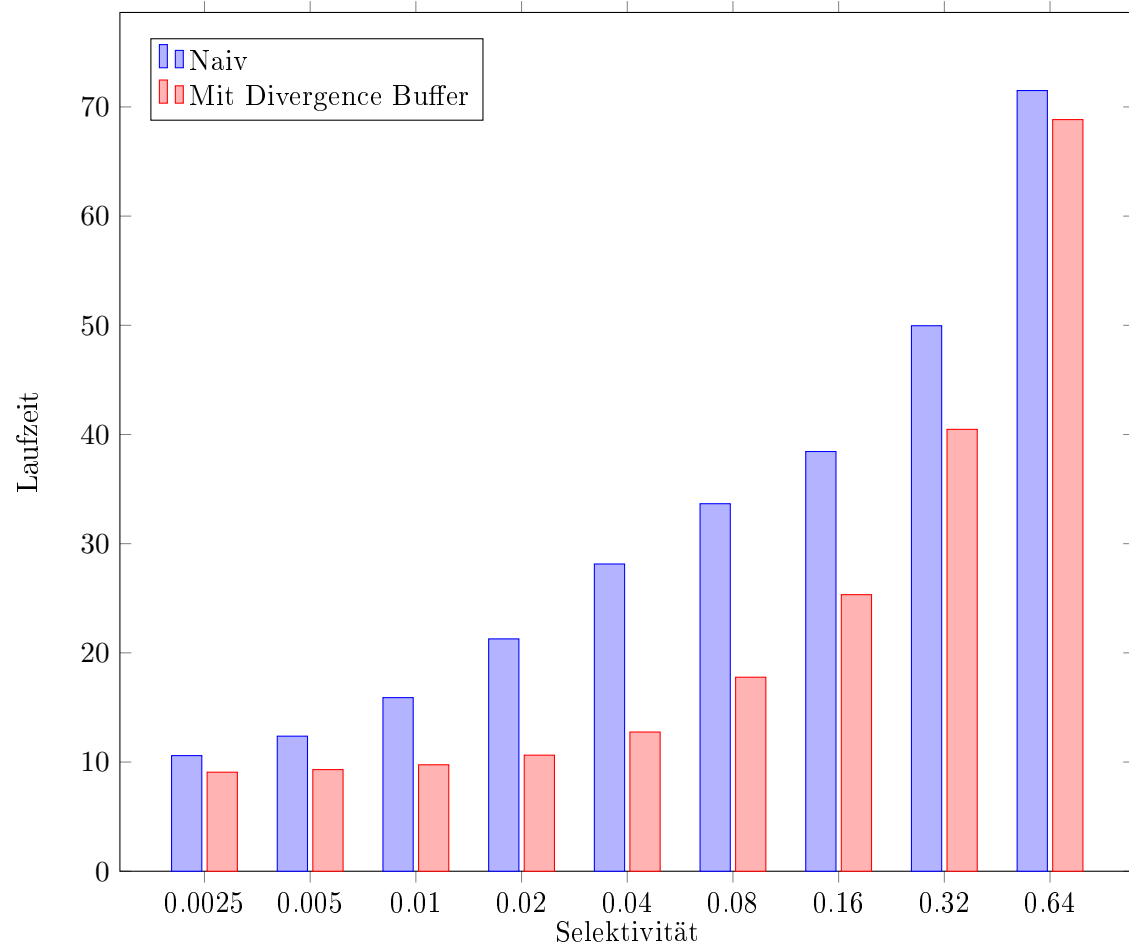
## 3.4 Präfixtest als alternativer Workload

## 3.5 Verwendete Workloads und deren Merkmale

## 3.6 Vorstellung der Messergebnisse







### 3.7 Diskussion der Ergebnisse

## Kapitel 4

# Paralleler Musterabgleich mit regulären Ausdrücken

- 4.1 Grundlagen zur Verarbeitung regulärer Ausdrücke
- 4.2 Naive Umsetzung
- 4.3 Umsetzung mit Lane Refill
- 4.4 Verwendete Workloads und deren Merkmale
- 4.5 Vorstellung der Messergebnisse
- 4.6 Diskussion der Ergebnisse





## Kapitel 5

### Ergebnis und Fazit



Anhang A

Weitere Informationen



# Abbildungsverzeichnis



Hiermit versichere ich, dass ich die vorliegende Arbeit selbstständig verfasst habe und keine anderen als die angegebenen Quellen und Hilfsmittel verwendet sowie Zitate kenntlich gemacht habe.

Dortmund, den 27. Januar 2019

Florian Lüdiger

