

# Übertragung eines an Mapreduce orientierten k-Means Algorithmus auf einen FPGA-Hardwarebeschleuniger für den Cloud-Einsatz

Studienarbeit

Martin Knöfel

Dresden, 12. Oktober 2017

## Ziele der Arbeit

- Literaturstudium zum Map-Reduce, Vivado High-Level-Synthese (VHLS) und k-Means Algorithmus
- Entwurf und Design des C++ Sourcecodes für VHLS
- Spezifizierung der Synthese mit VHLS Direktiven
- Implementierung und Auswertung

# Gliederung des Vortrags

1. Motivation
2. Literatur
3. Entwurfsprozess (HLS)
4. Auswertung
5. Zusammenfassung und Ausblick

1. Motivation

2. Literatur

3. Entwurfsprozess (HLS)

4. Auswertung

5. Zusammenfassung und Ausblick

## 1 Motivation

# k-Means Clustering

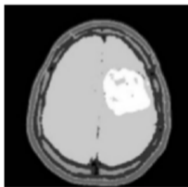
- übersichtliches Clusterverfahren
- großer Anwendungsbereich:
  - *Bildverarbeitung*
  - *Medizin*
  - *Statistik*
  - *Machine Learning*
- hoher/steigender Bedarf an Datensegmentierung

# 1 Motivation

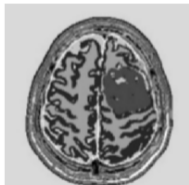
## k-Means Beispiel



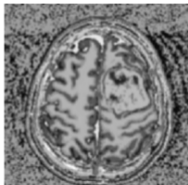
$k = 2$



$k = 6$



$k = 15$



$k = 40$

## 1 Motivation

# Cloud-Einsatz

- geringe Einstiegskosten für kleine Unternehmen
- spezielle Hardware (GPUs, FPGAs, ...)
- Skalierbarkeit der Ressourcen
- effiziente Nutzung der Hardware (zeitlich und räumlich)

# 1 Motivation

## High-Level-Synthese

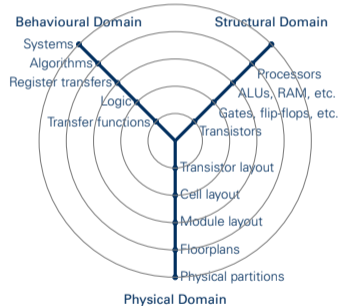


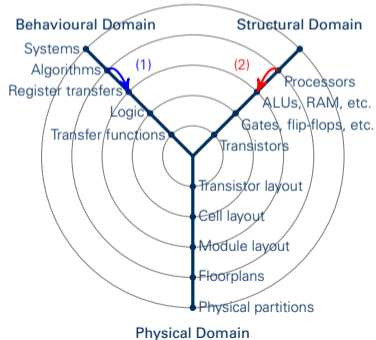
Abbildung: Gajski-Kuhn Y-Diagramm <sup>1</sup>

<sup>1</sup>Hemani, 2004: "Charting the EDA roadmap"



## 1 Motivation

# Vivado High-Level-Synthese



- (1) generiert VHDL/Verilog Code aus C/C++ Programm
- (2) Spezifikation der Hardwarestruktur über Direktiven(Pragmas)

## 1 Motivation

### IDEE:

- grundlegenden Struktur des k-Means Algorithmus
- Anpassung an den Bedarfsfall (Anwendung/Cloud)
- effiziente Implementierung durch die High-Level-Synthese

1. Motivation

2. Literatur

3. Entwurfsprozess (HLS)

4. Auswertung

5. Zusammenfassung und Ausblick

## 2 Literatur

# Map-Reduce

- Programmiermodell von Google <sup>2</sup>
- automatisierte Aufteilung von Daten
- zwei Schritte:

**MAP:** Verteilung der Daten auf die Verarbeitungsknoten

**REDUCE:** Zusammenfassen der Teilergebnisse

- k-Means → Abstandsberechnung, Zuordnung <sup>3</sup>

---

<sup>2</sup>Dean u. a., 2008: "MapReduce: simplified data processing on large clusters"

<sup>3</sup>Choi u. a., 2014: "Map-reduce processing of k-means algorithm with FPGA-accelerated computer cluster"

## 2 Literatur

# k-Means Allgemein

- Aufteilung von Daten in k-Cluster
- mittels Abstandsberechnung und Zuordnungsmethode
- terminierend und konvergiert zu einem lokalen Optimum <sup>4</sup>
- Standard-Algorithmus: Lloyd-Algorithmus

---

<sup>4</sup>Selim u. a., 1984: "K-means-type algorithms: A generalized convergence theorem and characterization of local optimality"

## 2 Literatur

## k-Means Abstandsberechnung

quadratische euklidische  
Distanz (Lloyd-Algorithmus):

$$d = \sum_{j=1}^k (a_j - b_j)^2$$

Manhattan-Distanz:

$$d = \sum_{j=1}^k |a_j - b_j|$$

Kosinus-Distanz:

$$d_{st} = 1 - \frac{x_s x_t'}{\sqrt{(x_s x_s')(x_t x_t')}}$$

Korrelations-Distanz:

$$d_{st} = 1 - \frac{(x_s - \bar{x}_s)(x_t - \bar{x}_t)}{\sqrt{(x_s - \bar{x}_s)(x_s - \bar{x}_s)'} \sqrt{(x_t - \bar{x}_t)(x_t - \bar{x}_t)'}}$$

5

<sup>5</sup>Bora u. a., 2014: "Effect of different distance measures on the performance of K-means algorithm: an experimental study in Matlab"

## 2 Literatur

# k-Means Pseudocode

```
while (Zentren != alte Zentren)
{
    for (alle Datenpunkte)
    {
        for (alle Zentren)
        {
            berechneAbstand(Datensatz , Zentrum);
        }
        kleinsten_Abstand ();
    }
    berechne_Zentren_neu ();
}
```

## 2 Literatur

### Verwandte Arbeiten

- Ein- und Ausgabe der Daten ist ein Flaschenhals
- Geschwindigkeit abhängig von der Menge der Cluster <sup>6 7</sup>
- High-Level-Synthese erzielt ungefähr gleiche Ergebnisse wie ein VHDL Entwurf
- „unroll“, „pipeline“ und „array\_partition“ Direktiven haben größten Einfluss <sup>8</sup>

---

<sup>6</sup>Choi u. a., 2014: “Map-reduce processing of k-means algorithm with FPGA-accelerated computer cluster”

<sup>7</sup>Lavenier, 2000: “FPGA implementation of the k-means clustering algorithm for hyperspectral images”

<sup>8</sup>Winterstein u. a., 2013: “High-level synthesis of dynamic data structures: A case study using Vivado HLS”



1. Motivation
2. Literatur
3. Entwurfsprozess (HLS)
4. Auswertung
5. Zusammenfassung und Ausblick

### 3 Entwurfsprozess (HLS)

## Anwendungsfall

- Bildsegmentierung im RGB Farbraum
  - 6 Cluster
  - dreidimensionaler Vektor
  - intuitive Kontrolle durch Visualisierung
  - 814 x 560 Pixel
- Reconfigurable Cloud Computing Framework (RC2F)<sup>9</sup> als Zielumgebung

---

<sup>9</sup>Knöfel u. a., 2015: “Computing framework for dynamic integration of reconfigurable resources in a cloud”

### 3 Entwurfsprozess (HLS)

## verfügbare Ressourcen (VIRTEX-7 485T)

Tabelle: verfügbare Ressourcen der RC2F Regionen

	total	Single	Dual	Triple	Quad	Quint	Hexa
<i>BRAM 18K</i>	2.060	100	200	300	400	500	600
<i>DSP48E</i>	2.800	340	660	980	1.320	1.660	1.940
<i>LUT</i>	303.600	30.800	60.400	90.000	120.800	151.600	188.400
<i>FF</i>	607.200	61.600	120.800	180.000	241.600	303.200	376.800

10

<sup>10</sup>Knodel u. a., 2017: "Virtualizing Reconfigurable Hardware to Provide Scalability in Cloud Architectures"

### 3 Entwurfsprozess (HLS)

## k-Means Funktion

- Parameter:
  - Vektorgroße
  - Bildgröße (Spalten/Zeilen)
  - Clusteranzahl
- Datentypen:
  - RGB Werte: 8Bit (unsigned char)
  - Berechnungen: 32Bit (long)

## 3 Entwurfsprozess (HLS)

# k-Means Funktion

```
void kMeans(data, result){
    loadData(data, result);
    do{
        copyOldCentroids();
        for(alles Pixel){
            for(alles Centroids){
                calcDistance(Pixel, Centroid);
            }
            accumulateMeans(assignment, pixel);
        }
        calcNewCentroids()
    }while(newCentroids!=oldCentroids);
}
```

## 3 Entwurfsprozess (HLS) Vorgehen

Simulate Design	<ul style="list-style-type: none"><li>• Validate The C function</li></ul>
Synthesize Design	<ul style="list-style-type: none"><li>• Baseline design</li></ul>
1: Initial Optimizations	<ul style="list-style-type: none"><li>• Define interfaces (and data packing)</li><li>• Define loop trip counts</li></ul>
2: Pipeline for Performance	<ul style="list-style-type: none"><li>• Pipeline and Dataflow</li></ul>
3: Optimize Structures for Performance	<ul style="list-style-type: none"><li>• Partition memories and ports</li><li>• Remove false dependencies</li></ul>
4: Reduce Latency	<ul style="list-style-type: none"><li>• Optionally specify latency requirements</li></ul>
5: Improve Area	<ul style="list-style-type: none"><li>• Optionally recover resources through sharing</li></ul>

11

### 3 Entwurfsprozess (HLS)

## Basisimplementierung

Tabelle: Ressourcenverbrauch

Name	BRAM_18K	DSP48	FF	LUT
<i>DSP</i>	-	1	-	-
<i>Expression</i>	-	-	0	642
<i>FIFO</i>	-	-	-	-
<i>Instance</i>	-	-	752	864
<i>Memory</i>	1.026	-	160	12
<i>Multiplexer</i>	-	-	-	334
<i>Register</i>	-	-	495	-
<i>Total</i>	1.026	1	1407	1852
<i>Available</i>	2.060	2.800	607.200	303.600
<i>Utilization (%)</i>	49	~0	~0	~0

### 3 Entwurfsprozess (HLS)

## Basisimplementierung

Tabelle: Latenzen in Takten

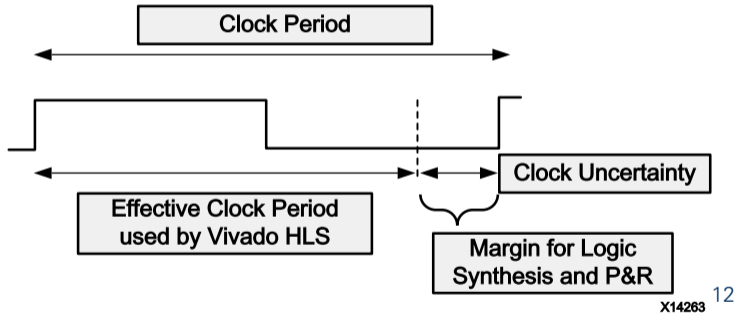
Latency		Interval	
<i>min</i>	<i>max</i>	<i>min</i>	<i>max</i>
30.085.724	30.086.912	30.085.725	30.086.913

Tabelle: Timing Basisimplementierung

	Target	Estimated	Uncertainty
Clock in ns	10,00	8,55	1,25



### 3 Entwurfsprozess (HLS) Konzept der Unsicherheit



## 3 Entwurfsprozess (HLS)

# Schleifenübersicht

Tabelle: Übersicht der Schleifen

Loop Name	Latency		Iteration Latency	Trip Count	Pipelined
	<i>min</i>	<i>max</i>			
- <i>load11</i>	3.646.768	3.646.768	8	455.846	no
+ <i>load_12</i>	6	6	2	3	no
+ <i>load_13</i>	6	6	2	3	no
- <i>dowhileBody</i>	26.438.954	26.440.142	26.438.955 ~ 26.440.143	0 ~ 1	no
+ <i>copyCEN11</i>	48	48	8	6	no
++ <i>copyCEN12</i>	6	6	2	3	no
+ <i>calcPixel</i>	26.438.720	26.438.720	58	455.840	no
++ <i>calcPCEN</i>	48	48	8	6	no
+++ <i>calcD</i>	6	6	2	3	no
++ <i>saveCenSum</i>	6	6	2	3	no
+ <i>dowhileBody.3</i>	48	48	8	6	no
++ <i>cenOut</i>	6	6	2	3	no
+ <i>compare11</i>	48	48	8	6	no
++ <i>compare12</i>	6	6	2	3	no

- = 1. Ebene

+ = 2. Ebene

++ = 3. Ebene

### 3 Entwurfsprozess (HLS)

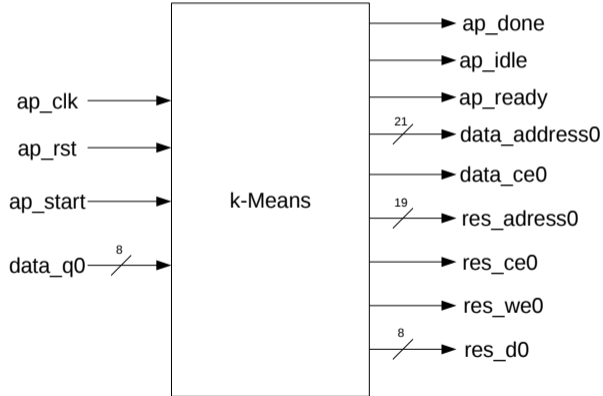
## Schleifenübersicht

Tabelle: Übersicht der Schleifen

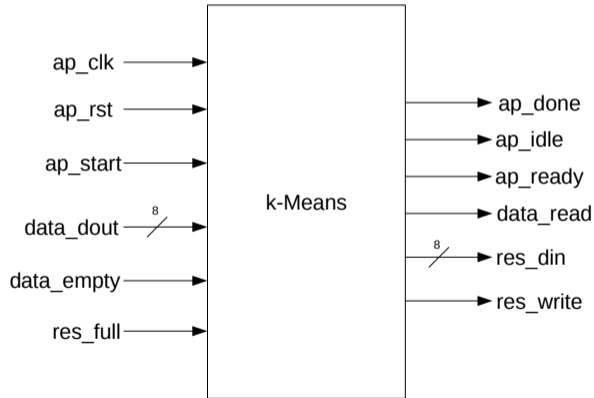
Loop Name	Latency		Iteration Latency	Trip Count	Pipelined
	<i>min</i>	<i>max</i>			
- <i>load11</i>	3.646.768	3.646.768	8	455.846	no
+ <i>load_12</i>	6	6	2	3	no
+ <i>load_13</i>	6	6	2	3	no
- <i>dowhileBody</i>	26.438.954	26.440.142	26.438.955 ~ 26.440.143	0 ~ 1	no

### 3 Entwurfsprozess (HLS)

## Basisinterface



### 3 Entwurfsprozess (HLS) Interface angepasst



### 3 Entwurfsprozess (HLS)

## Direktive: „PIPELINE“

- für Schleifen und Funktionen
- Verwendung: von Innen nach Außen
- nutzt automatisch die „unroll“ Direktive
- Optionen:
  - *Iteration Interval (II)*
  - *enable loop rewinding*
  - *enable flushing*
  - *disable loop pipelining*
- globales Flag möglich

### 3 Entwurfsprozess (HLS)

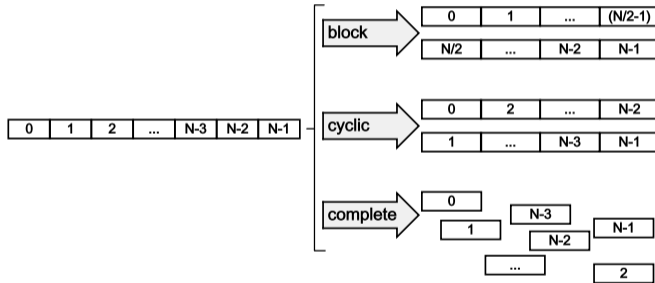
## Direktive: „PIPELINE“

Tabelle: Vergleich des Ressourcenverbrauchs

	<i>FIFO</i>	<i>Pipeline</i>	<i>Inline 10ns</i>	<i>Inline 20ns</i>
BRAM_18K	1.026	773	774	774
DSP48	1	3	18	18
FF	1.407	1.577	9.323	9.177
LUT	1.852	2.007	9.475	9.473

### 3 Entwurfsprozess (HLS)

## Direktive: „ARRAY\_PARTITION“



13



### 3 Entwurfsprozess (HLS)

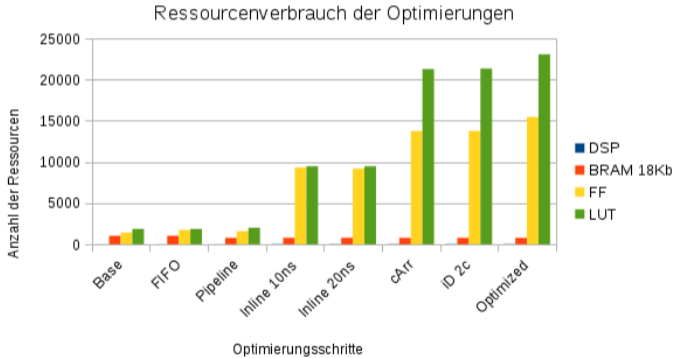
## Direktive(GUI): „config\_core“

- IP Cores werden ohne Pipeline verwendet
- bei bestimmten IP Cores können mit Pipelines initialisiert werden
- z.B. der DSP48 Kern (mit einer Latenz von 4 Takten)
- wirkt sich sehr positiv auf das Scheduling der Synthese aus

1. Motivation
2. Literatur
3. Entwurfsprozess (HLS)
4. Auswertung
5. Zusammenfassung und Ausblick

## 4 Auswertung

# Optimierungsschritte:



## 4 Auswertung Simulation

- RTL/Co Simulation
- erzeugt automatisch eine VHDL Testbench aus der C++ Testbench
- Datenausgaben werden separat gespeichert
- Waveform Ausgabe für Ein- und Ausgänge <sup>14</sup>

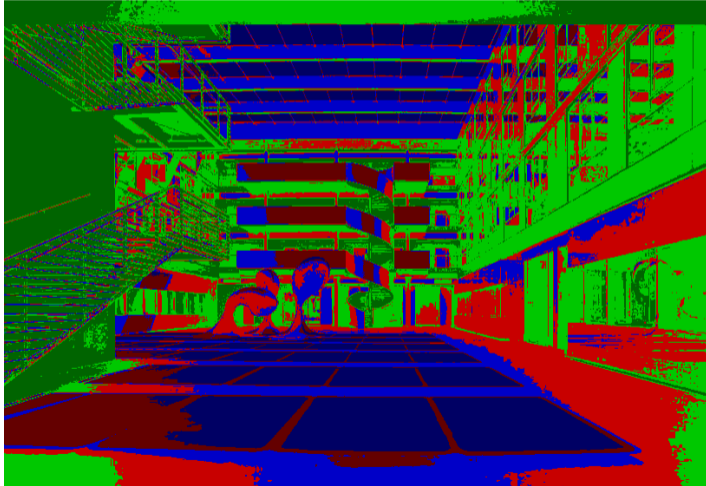
---

<sup>14</sup>, 2017: [Vivado Design Suite User Guide, High Level Synthesis](#)

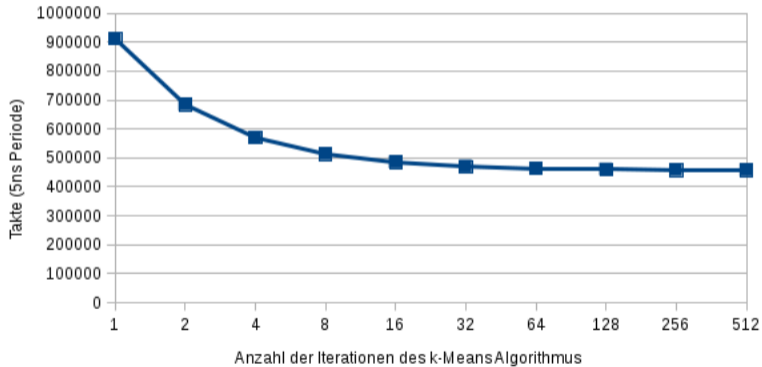
## Abbildung: original Testbild



## Abbildung: Ergebnis der Simulation



### Latenz pro Durchlauf



## 4 Auswertung

# Schleifenübersicht

Tabelle: Übersicht der Schleifen

<i>Loop Name</i>	<i>Latency</i>		<i>Iteration Latency</i>	<i>Initiation Interval</i>		<i>Trip Count</i>	<i>Pipelined</i>
	<i>min</i>	<i>max</i>		<i>achieved</i>	<i>target</i>		
- <i>load11</i>	455.847	455.847	2	1	1	455.846	yes
- <i>dowhileBody</i>	455.990	4.559.909	455.991	-	-	1 ~ 10	no
+ <i>copyCEN11</i>	6	6	1	1	1	6	yes
+ <i>calcPixel</i>	455.851	455.851	13	1	1	455.840	yes
+ <i>newCEN11_newCEN12</i>	87	87	71	1	1	18	yes
+ <i>L_cenOut</i>	18	18	2	1	1	18	yes
+ <i>compare11_compare12</i>	18	18	2	1	1	18	yes



## 4 Auswertung

# Clusteranzahl

Tabelle: Anpassung der Clusteranzahl

		<i>K=2</i>	<i>K=3</i>	<i>K=4</i>	<i>K=5</i>	<i>K=6</i>	<i>K=7</i>	<i>K=8</i>	<i>K=16</i>	<i>K=32</i>	<i>K=64</i>
BRAM_18K		768	768	768	768	768	768	768	768	768	2.060
DSP48		6	9	12	15	18	21	24	48	96	1
FF		10.915	11.106	12.052	14.277	15.461	16.399	17.054	23.240	37.851	80.739
LUT		11.887	14.174	16.156	19.454	23.082	25.990	29.520	38.497	85.393	<b>401.834</b>
Clock	<i>Target</i>	5,00	5,00	5,00	5,00	5,00	5,00	5,00	5,00	5,00	5,00
	<i>Estimated</i>	4,87	<b>5,39</b>	<b>5,39</b>	<b>5,57</b>	4,36	4,26	4,44	4,26	4,26	4,33
Latenz	<i>min</i>	911.793	911.805	911.816	911.829	911.840	911.853	911.864	911.956	912.140	183.248.529
	<i>max</i>	5.015.325	5.015.436	5.015.537	5.015.649	5.015.759	5.015.871	5.015.972	5.016.820	5.018.516	1.824.279.000

## 4 Auswertung

# Clusteranzahl

Tabelle: Anpassung der Clusteranzahl (Ausschnitt)

		<i>K=2</i>	<i>K=3</i>	<i>K=4</i>	<i>K=5</i>
BRAM_18K		768	768	768	768
DSP48		6	9	12	15
FF		10.915	11.106	12.052	14.277
LUT		11.887	14.174	16.156	19.454
Clock	<i>Target</i>	5,00	5,00	5,00	5,00
	<i>Estimated</i>	4,87	<b>5,39</b>	<b>5,39</b>	<b>5,57</b>
Latenz	<i>min</i>	911.793	911.805	911.816	911.829
	<i>max</i>	5.015.325	5.015.436	5.015.537	5.015.649

## 4 Auswertung

# Ergebnis

- Bildbeispiel: 80 Iterationen
- Taktperiode: 5 ns (200 MHz)
- Gesamtlatenz: 36.935.130
- = 184,6 ms
- C++ Sourcecode des k-Means Algorithmus
- i5-4690 (3,5GHz)
- 1 Kern
- = 3,704 s

→ FPGA Implementierung 20x schneller als eigene C++ Implementierung

1. Motivation
2. Literatur
3. Entwurfsprozess (HLS)
4. Auswertung
5. Zusammenfassung und Ausblick

## 5 Zusammenfassung und Ausblick

# Vivado HLS

### Vorteile:

- effizientes Pipelining
- schnelles Anpassen der Hardware
- bessere Lesbarkeit

### Nachteile:

- keine vollständige High-Level-Synthese (Direktiven, Optionen)
- Fachkenntnisse für Hardware und Synthese notwendig
- Details manchmal schwer umsetzbar
- keine Unterstützung für Sourcecode Änderungen

## 5 Zusammenfassung und Ausblick

### Ausblick

- Testen einer modularen Variante
- weitere Untersuchungen zum Verhalten
- Regeln für Synthese aufstellen bzgl. Anwendungen und Ressourcen
- Entwicklung einer Software für die automatisierte Synthese des k-Means Algorithmus
- Neuronale Netze für automatische High-Level-Synthese