

KOMPRESSION

WIE HILFT DATENKOMPRESSION BEI DER VERARBEITUNG VON BIG DATA?

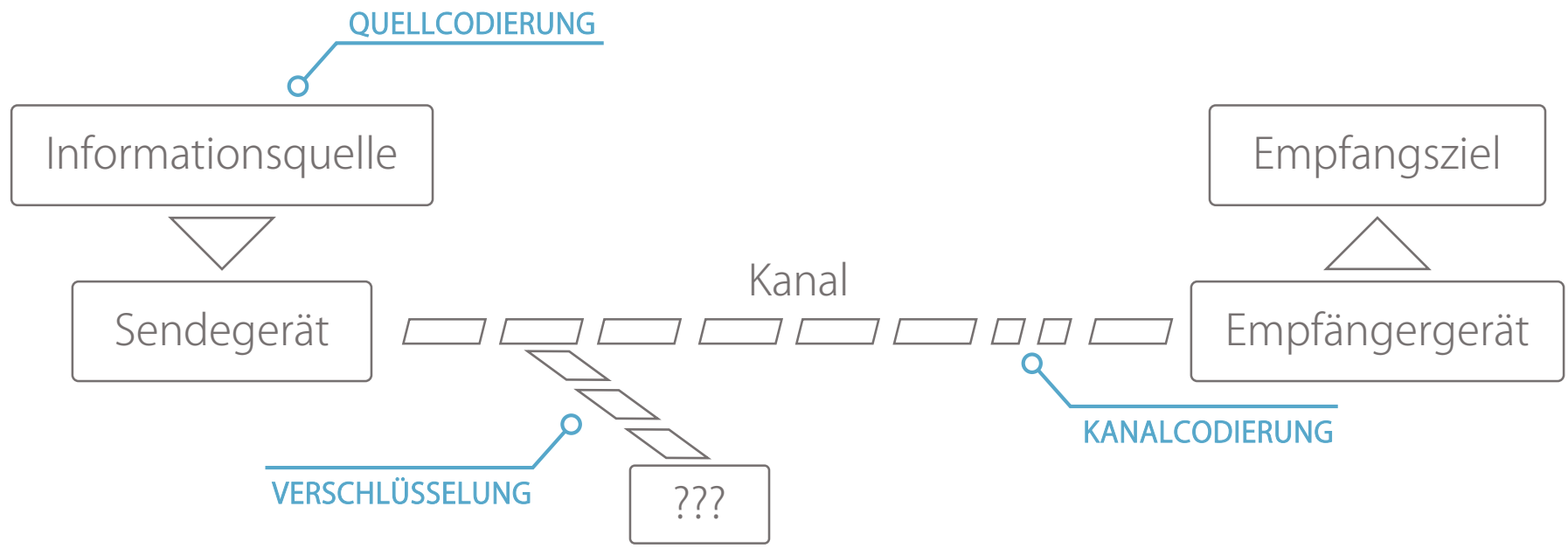
MARVIN AHMADI-MOGHADDAM | PROSEMINAR EFFIZIENTE PROGRAMMIERUNG

UNIVERSITÄT HAMBURG | 07.06.2018

GLIEDERUNG

- I WAS IST KOMPRESSION IN DER NACHRICHTENTHEORIE?
- II ARTEN DER KOMPRESSION UND KOMPRESSIONSVERFAHREN
- III KOMPRESSION UND BIG DATA

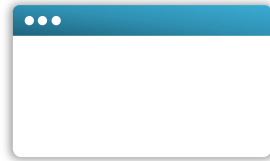
KOMMUNIKATIONSMODELL NACH SHANNON



VERLUSTFREIE & VERLUSTBEHAFTETE KOMPRESSION

VERLUSTFREIE KOMPRESSION

EINDEUTIG KORREKTE WIEDERHERSTELLUNG



PROGRAMME



ÜBERTRAGUNG ÜBER
DAS INTERNET

VERLUSTFREIE KOMPRESSION

CODIERUNGSVERFAHREN

PRINZIP: INFORMATIONEN ANDERS REPRÄSENTIEREN UND REDUNDANZ VERMINDERN

REPRÄSENTATION:

5 / V

01 / [NIEDRIGE SPANNUNG] [HOHE SPANNUNG]

CODIERUNGSVERFAHREN

BLOCKCODE: FESTE LÄNGE

0 0 0 0 1 1 0 0 0 0 0 1 1 1 1 1 1 0 1 1 1 0 0 0 0 0 0 1

0 0 0 0	0 0 0 1	1 0 1 1
1 1 0 0	1 1 1 1	1 0 0 0

28 ZEICHEN IM BITSTRING

4 ZEICHEN PRO WORT

7 WÖRTER IM BITSTRING

VERLUSTFREIE KOMPRESSION

CODIERUNGSVERFAHREN

PRÄFIXCODE: VARIABLE LÄNGE

0 0 0 0 1 1 0 0 0 0 0 1 1 1 1 1 1 0 1 1 1 0 0 0 0 0 0 1

0 0 0 0	0 0 1	1 0
0 0 0 1	1 1	0 1

28 ZEICHEN IM BITSTRING

2 BIS 4 ZEICHEN PRO WORT

11 WÖRTER IM BITSTRING

PRÄFIX-EIGENSCHAFT: EINDEUTIG DEKODIERBAR

VERLUSTFREIE KOMPRESSION

CODIERUNGSVERFAHREN

LAUFLÄNGENCODIERUNG: REDUNDANZ VERMINDERN

0000	1100	0001	1111	1011	1000	0001
40	21 50	61		1031	60	11

28 ZEICHEN IM BITSTRING

16 ZEICHEN LAUFLÄNGENKODIERT

VERLUSTFREIE KOMPRESSION

CODIERUNGSVERFAHREN

ENTROPIEKODIERUNG: HÄUFIGE WÖRTER KURZ KODIERT

a) SHANNON-FANO

b) HUFFMAN

c) ARITHMETISCH

CODIERUNGSVERFAHREN

HUFFMAN:

A: 50%

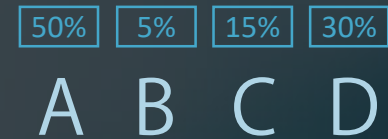
B: 5%

C: 15%

D: 30%

-
- 1.) WÖRTER AUFSTEIGEND NOTIEREN
 - 2.) WÖRTER DER GERINGSTEN WAHRSCHEINLICHKEIT ZU NEUEM WORT ZUSAMMENFASSEN
 - 3.) WIEDERHOLEN BIS 2 WÖRTER ÜBRIG BLEIBEN

=> KLEINSTMÖGLICHE MITTLERE CODEWORTLÄNGE



CODIERUNGSVERFAHREN

HUFFMAN:

A: 50%

B: 5%

C: 15%

D: 30%

-
- 1.) WÖRTER AUFSTEIGEND NOTIEREN
 - 2.) WÖRTER DER GERINGSTEN WAHRSCHEINLICHKEIT ZU NEUEM WORT ZUSAMMENFASSEN
 - 3.) WIEDERHOLEN BIS 2 WÖRTER ÜBRIG BLEIBEN

=> KLEINSTMÖGLICHE MITTLERE CODEWORTLÄNGE



CODIERUNGSVERFAHREN

HUFFMAN:

A: 50%

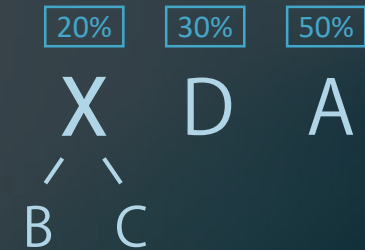
B: 5%

C: 15%

D: 30%

-
- 1.) WÖRTER AUFSTEIGEND NOTIEREN
 - 2.) WÖRTER DER GERINGSTEN WAHRSCHEINLICHKEIT ZU NEUEM WORT ZUSAMMENFASSEN
 - 3.) WIEDERHOLEN BIS 2 WÖRTER ÜBRIG BLEIBEN

=> KLEINSTMÖGLICHE MITTLERE CODEWORTLÄNGE



CODIERUNGSVERFAHREN

HUFFMAN:

A: 50%

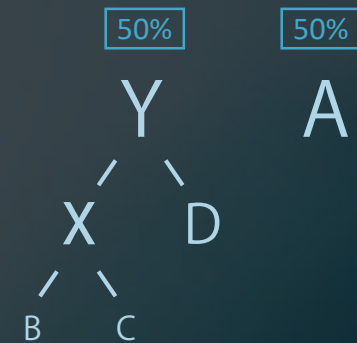
B: 5%

C: 15%

D: 30%

- 1.) WÖRTER AUFSTEIGEND NOTIEREN
- 2.) WÖRTER DER GERINGSTEN WAHRSCHEINLICHKEIT ZU NEUEM WORT ZUSAMMENFASSEN
- 3.) WIEDERHOLEN BIS 2 WÖRTER ÜBRIG BLEIBEN

=> KLEINSTMÖGLICHE MITTLERE CODEWORTLÄNGE



CODIERUNGSVERFAHREN

HUFFMAN:

A: 50%

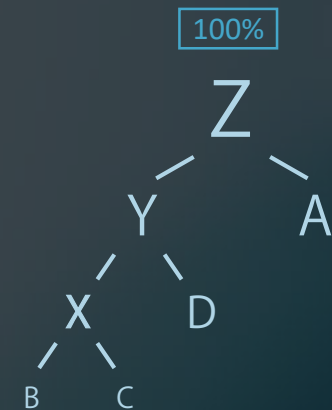
B: 5%

C: 15%

D: 30%

-
- 1.) WÖRTER AUFSTEIGEND NOTIEREN
 - 2.) WÖRTER DER GERINGSTEN WAHRSCHEINLICHKEIT ZU NEUEM WORT ZUSAMMENFASSEN
 - 3.) WIEDERHOLEN BIS 2 WÖRTER ÜBRIG BLEIBEN

=> KLEINSTMÖGLICHE MITTLERE CODEWORTLÄNGE



CODIERUNGSVERFAHREN

HUFFMAN:

A: 50%

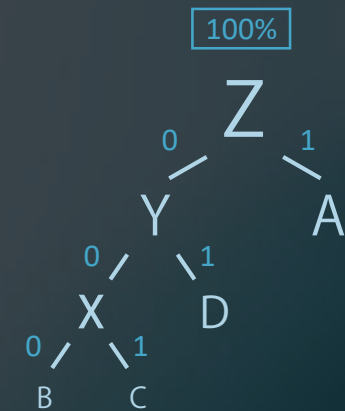
B: 5%

C: 15%

D: 30%

-
- 1.) WÖRTER AUFSTEIGEND NOTIEREN
 - 2.) WÖRTER DER GERINGSTEN WAHRSCHEINLICHKEIT ZU NEUEM WORT ZUSAMMENFASSEN
 - 3.) WIEDERHOLEN BIS 2 WÖRTER ÜBRIG BLEIBEN

=> KLEINSTMÖGLICHE MITTLERE CODEWORTLÄNGE



VERLUSTFREIE KOMPRESSION

KOMPRESSIONSVERFAHREN

DICTIONARY-BASED

- TEILE EINES STRINGS TAUCHEN IM TEXT MEHRMALS AUF
- DYNAMISCH VS. STATISCH

LZ77

ABRAHAM LEMPEL, JACOB ZIV
1977

VERLUSTFREIE KOMPRESSION

KOMPRESSIONSVERFAHREN

SLIDING WINDOW (LZ77)



VERLUSTFREIE KOMPRESSION

KOMPRESSIONSVERFAHREN

SLIDING WINDOW (LZ77)



VERLUSTFREIE KOMPRESSION

KOMPRESSIONSVERFAHREN

SLIDING WINDOW (LZ77)



VERLUSTFREIE KOMPRESSION

KOMPRESSIONSVERFAHREN

SLIDING WINDOW (LZ77)



(5 | |)
OFFSET LÄNGE ZEICHEN

VERLUSTFREIE KOMPRESSION

KOMPRESSIONSVERFAHREN

SLIDING WINDOW (LZ77)



(5 | 3 | L)
OFFSET LÄNGE ZEICHEN

VERLUSTFREIE KOMPRESSION

KOMPRESSIONSVERFAHREN

SLIDING WINDOW (LZ77)



VERLUSTFREIE KOMPRESSION

KOMPRESSIONSVERFAHREN

SLIDING WINDOW (LZ77)



(0 | 0 | P)
OFFSET LÄNGE ZEICHEN

VERLUSTFREIE KOMPRESSION

KOMPRESSIONSVERFAHREN

BEISPIELE AUS DEM ALLTAG:



ZIP



PNG

VERLUSTBEHAFTETE KOMPRESSION

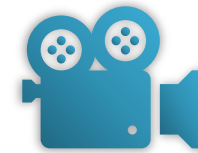
PRINZIP: IRRELEVANTE DATEN FILTERN



BILDKOMPRESSION



AUDIOKOMPRESSION



VIDEOKOMPRESSION

VERLUSTBEHAFTETE KOMPRESSION

BILDKOMPRESSION

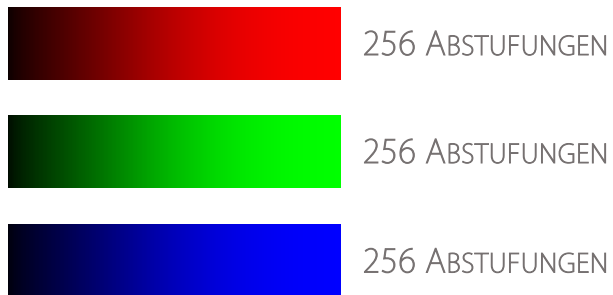
8 ABSTUFUNGEN



16 ABSTUFUNGEN



8-BIT RGB-MODELL



16.777.216 KOMBINATIONEN

MENSCHLICHES AUGE



20.000.000
VERSCHIEDENE FARBEN

VERLUSTBEHAFTETE KOMPRESSION

BILDKOMPRESSION

PRINZIP DER NATÜRLICHEN BILDKOMPRESSION:

WENN WIR EINEN PIXEL IN EINEM BILD ZUFÄLLIG WÄHLEN, IST DIE WAHRSCHEINLICHKEIT HOCH, DASS BENACHBARTE PIXEL DIE SELBE FARBE ODER ÄHNLICHE FARBEN HABEN

[SAL]

VERLUSTBEHAFTETE KOMPRESSION

BILDKOMPRESSION

NATÜRLICHE BILDER

Y: LUMINANCE

CB, CR: COLOR

-> MÖGLICHEST GERINGE VERLUSTE

-> VERLUSTE DÜRFEN HÖHER SEIN

EINDEUTIGE FARBABSTUFUNGEN

ÄHNLICHE BEREICHE IDENTIFIZIEREN UND DARAUFG VERWEISEN

DATENKOMPRESSION IN BIG DATA

6 Vs VON BIG DATA

VOLUME:

GRÖßE, DIE FÜR VORLIEGENDE
COMPUTER NICHT MEHR
EINFACH ZU VERARBEITEN SIND

VELOCITY:

MENGE AN STRÖMEN,
DIE ZU EINEM BIG DATA
STROM FLIEßEN

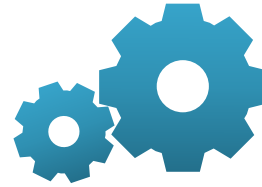
VARIETY:

VERSCHIEDENE
DATENQUELLEN IN
VERSCHIEDENEN FORMATEN

KOMPLEXITÄT VON BIG DATA



1.) KOMPLEXITÄT DER **DATEN**



2.) KOMPLEXITÄT DER **VERARBEITUNG**



3.) KOMPLEXITÄT DER **SYSTEME**

=> **AUSSORTIEREN** UND **ZUSAMMENFASSEN**

VERMINDERUNG DER DATEN

VOLUME:

GRÖßE, DIE FÜR VORLIEGENDE
COMPUTER NICHT MEHR
EINFACH ZU VERARBEITEN SIND

VELOCITY:

MENGE AN STRÖMEN,
DIE ZU EINEM BIG DATA
STROM FLIEßEN

VARIETY:

VERSCHIEDENE DATENQUELLEN
IN VERSCHIEDENEN FORMATEN

VERMINDERUNG DER
DATENMENGE

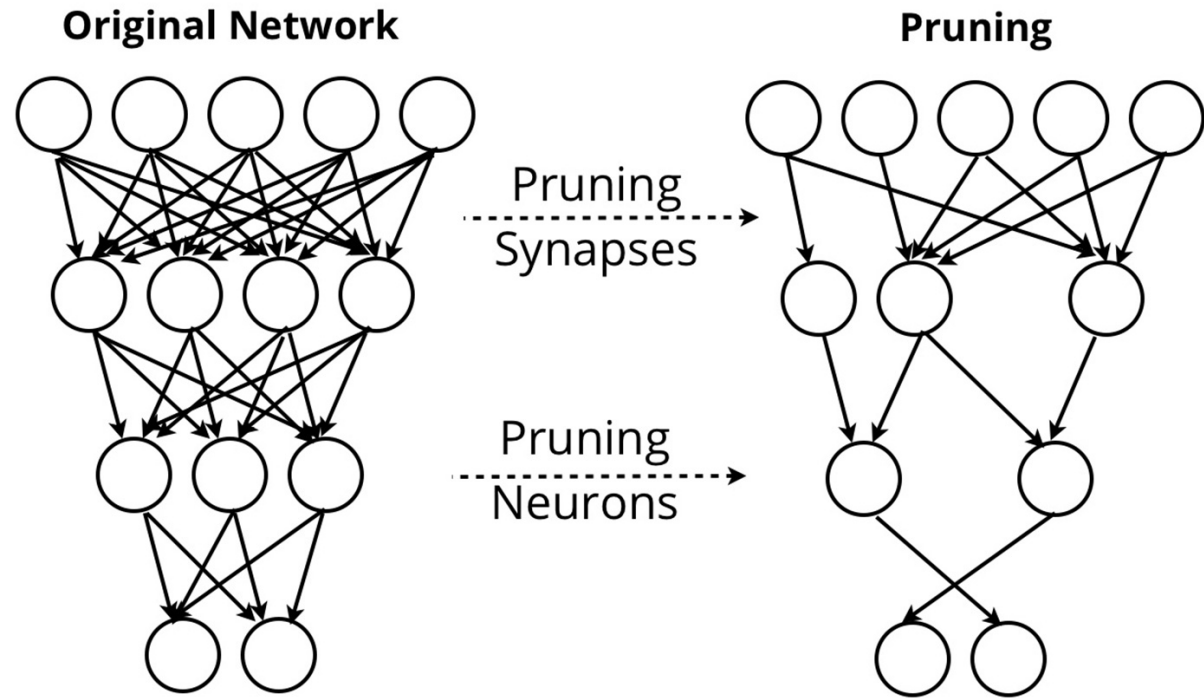
ZUSAMMENFASSEN
UND SORTIEREN ZUR
VERMINDERUNG DER
REDUNDANZ

ZUSAMMENFASSEN UM
HETEROGENITÄT ZU
VERMINDERN

=> PRINZIP DER **MINIMIERUNG VON REDUNDANZ**

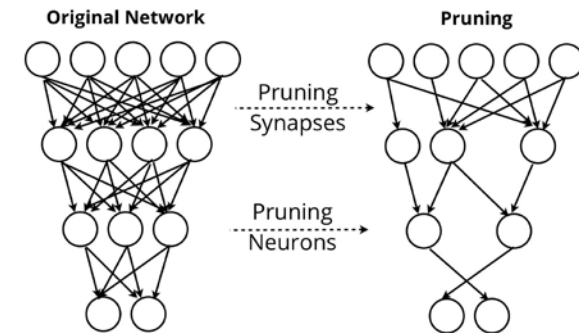
DEEP LEARNING BEISPIEL

- 1.) Entfernen der überflüssigen Verbindungen
- 2.) Quantisierung und Gewichtung
- 3.) Huffman-Kodierung



DEEP LEARNING BEISPIEL

- 1.) Entfernen der überflüssigen Verbindungen
- 2.) Quantisierung und Gewichtung
- 3.) Huffman-Kodierung



DATEIGRÖÖE VERMINDERN
35 BIS 50-MAL

ECHTZEITÜBERTRAGUNG MÖGLICH
HÄUFIGES AKTUALISIEREN AUF MOBIELN GERÄTEN

FAZIT

KOMPRESSION

WIE HILFT DATENKOMPRESSION BEI DER VERARBEITUNG VON BIG DATA?

QUELLEN

Lorica (2016) **Compressed representations in the age of big data** O'REILLY.

Abgerufen am 05.06.2018 von: <https://www.oreilly.com/ideas/compressed-representations-in-the-age-of-big-data>

ur Rehman, Liew, Abbas, Jayaraman, Wah, Khan (2016) **Big Data Reduction Methods: A Survey** Springer Berlin Heidelberg.

Abgerufen am 05.06.2018 von <https://link.springer.com/article/10.1007/s41019-016-0022-0>

The History of Data Compression (11.11.2014) **Commander**

Abgerufen am 05.06.2018 von <http://blog.commander.com/the-history-of-data-compression/>

Shannon (1948) **A Mathematical Theory of Communication** The Bell System Technical Journal.

Abgerufen am 05.06.2018 von <https://link.springer.com/article/10.1007/s41019-016-0022-0>

The History of Data Compression (11.11.2014) **Commander**

Abgerufen am 05.06.2018 von <http://blog.commander.com/the-history-of-data-compression/>

History of Lossless Data Compression Algorithms (20.08.2014) **Engineering and Technology History Wiki**

Abgerufen am 05.06.2018 von http://ethw.org/History_of_Lossless_Data_Compression_Algorithms

Wie viele Farben sieht der Mensch? (15.05.2014) **Schweizer Radio und Fernsehen**

Abgerufen am 05.06.2018 von <https://www.srf.ch/sendungen/einstein/fuenfmalklug/wie-viele-farben-sieht-der-mensch>

Datenkompression (02.05.2018) **Wikipedia**

Abgerufen am 05.06.2018 von <https://de.wikipedia.org/wiki/Datenkompression>

Salomon (2008) **A Concise Introduction to Data Compression** Springer-Verlag London