



Hochschule Darmstadt
- Fachbereich Informatik -

Grundlagen der Videokompression

Seminararbeit im Kurs
Wissenschaftliches Arbeiten in der Informatik I

vorgelegt von
Justin Böhm und Matthias Greune

Referentin: <Name>

Ausgabedatum: <Datum>

Abgabedatum: <Datum>

Erklärung

Ich versichere hiermit, dass ich die vorliegende Arbeit selbständig verfasst und keine anderen als die im Literaturverzeichnis angegebenen Quellen benutzt habe. Alle Stellen, die wörtlich oder sinngemäß aus veröffentlichten oder noch nicht veröffentlichten Quellen entnommen sind, sind als solche kenntlich gemacht. Die Zeichnungen oder Abbildungen in dieser Arbeit sind von mir selbst erstellt worden oder mit einem entsprechenden Quellennachweis versehen. Diese Arbeit ist in gleicher oder ähnlicher Form noch bei keiner anderen Prüfungsbehörde eingereicht worden.

<Name>

<Ort>, den 5. Dezember 2016

Abstrakt

Diese Arbeit gibt einen Überblick über die grundlegenden Methoden von Videokompressionsverfahren. Hierfür werden, sich am Encoding-Prozess von MPEG-1 orientierend, zunächst Arten der Irrelevanzreduktion, anschließend die wichtigsten Ansätze der Redundanzreduktion vorgestellt und anhand von Beispielcode erläutert. In eigenen Tests wurden unter Anwendung der vorgestellten Methoden zur verlustbehafteten und partieller Anwendung von verlustfreien Kompressionsalgorithmen Kompressionsraten mit einer Ratio bis zu 1:60 erreicht. Diese Arbeit zeigt somit, wie mittels weniger Grundlagen bereits vergleichsweise hohe Einsparungen im Speicherverbrauch von Videos erreicht werden können. Insbesondere mit Blick auf die stetig steigenden Forderungen nach höheren Framerates, und besserer Auflösung wird deutlich, welche hohe Relevanz das Thema Videokompression auch in Zukunft haben wird.

Inhaltsverzeichnis

Erklärung	ii
Abstrakt	iii
Abbildungsverzeichnis	v
1. Einleitung	1
2. Irrelevanzreduktion	2
2.1. Chroma Subsampling	2
2.2. Diskrete Kosinus Transformation	3
2.3. Quantisierung	4
3. Redundanzreduktion	6
3.1. Entropiecodierung	6
3.2. Motion Compensation	6
3.2.1. Frames	7
3.2.2. Makroblocks	8
3.2.3. Motion Estimation and Compensation	8
4. Ausblick	9
5. Zusammenfassung	10
A. Weitere Abbildungen und Tabellen	15
Literatur	17

Abbildungsverzeichnis

2.1. Mittels DCT gut komprimierbarer 8x8 Pixelblock	4
A.1. Artefakte durch Chroma Subsampling	15
A.2. Ergebnis der Quantisierung mit verschiedenen Quantisierungsfaktoren . . .	16

1. Einleitung

Videos sind seit der Entwicklung des Fernsehers zum Massenmedium kaum noch aus dem alltäglichen Leben wegzudenken. Seit dem Aufstieg des Internets als zentrales Kommunikationsmedium haben sich allerdings die Anforderungen an geeignete Speichertechniken von Videos drastisch verändert. Die heutigen Abspielgeräte haben noch immer begrenzten Speicherplatz und sind häufig nur mit schmalbandigen Internetanbindungen ausgestattet. Die Auflösung der Videos ist hingegen stark gestiegen. Um diese Ansprüche zu adressieren wurden Kompressionsalgorithmen entwickelt, die eine effiziente Speicherung speziell für bewegte Bilder ermöglichen. Die resultierenden Probleme aus dieser Art der Speicherung, wie Bildartefakte, sind heutigen Nutzern wohlbekannt. Die eigentliche Funktionsweise von Videokompression bleibt aber oft unbemerkt. TBC

2. Irrelevanzreduktion

Die rohe Aufnahme eines Bildes bietet eine Fülle an Informationen. Mit Blick auf die Eigenschaften des menschlichen Sehsinns lässt sich hierbei allerdings feststellen, dass einige Informationen relevanter für das Erkennen eines Bildes sind, als andere. Die Irrelevanzreduktion beschäftigt sich mit der Trennung und Reduzierung von weniger wichtigen Informationen und bietet damit Methoden zur verlustbehafteten Datenkompression an.

Bei der Videokompression werden im wesentlichen zwei Eigenschaften zur Reduktion von Daten ausgenutzt. Zum einen nimmt das Auge Varianzen in der Helligkeit (Luminanz) stärker wahr, als Änderungen im Farbton (Chrominanz). Zum Anderen ist das Auge besser in der Lage niedrige Ortsfrequenzen zu erkennen, als hohe - erkennt also grobe Strukturen eher als feinere. Diese Eigenschaften können nun ausgenutzt werden, um einen guten Kompromiss aus akzeptabler Bildqualität und guter Datenreduktion zu finden [Akr14].

2.1. Chroma Subsampling

Das Chroma Subsampling nutzt den Umstand aus, dass Helligkeitsvarianzen besser wahrgenommen werden, als Farbvarianzen. Zumeist liegen die Bildinformationen im Ausgangsformat jedoch im RGB Farbmodell vor, wobei hier die Helligkeitswerte in jeden Kanal eingehen. Um nun aber die Chrominanz bei gleichbleibender Auflösung der Luminanz zu reduzieren wird eine getrennte Darstellung dieser Informationen benötigt. Hierfür wird im MPEG-1 Standard die $YC_B C_R$ Darstellung verwendet, wobei das Y für die Luminanz steht und in C_B und C_R die Farbwerte codiert werden. Die Umrechnung lässt sich mittels folgender Formeln realisieren [Int95]:

$$Y = 0.299 \cdot R + 0.587 \cdot G + 0.114 \cdot B$$

$$U = (B - Y) \cdot 0.493$$

$$V = (R - Y) \cdot 0.877$$

Nun kann das eigentliche Subsampling stattfinden, welches bei MPEG-1 bei einer Auflösung von 4:2:0 realisiert wird. Die erste Zahl gibt hierbei die horizontale Abtastrate der Luminanz an. Die zweite Zahl steht für die horizontale Abtastrate der C_B und C_R Kanäle in Relation zum ersten Wert. Die dritte Zahl gibt die vertikale Samplingrate an, wobei

diese entweder 2 oder 0 betragen kann, also entweder kein vertikales Subsampling, oder vertikales Subsampling von 2:1 stattfindet. Für den Fall von 4:2:0 Subsampling bedeutet dies, dass jeweils 2x2 Bildpunkte des C_B und C_R Kanals auf einen Bildpunkt in der Ergebnismenge abgebildet werden. Hiermit wird also die Auflösung des C_B und C_R Kanals halbiert, was zu einer Datenreduktion von insgesamt 50% führt. [Poy]

Das Chroma Subsampling bietet somit eine gute Möglichkeit der Kompression, die allerdings nicht verlustfrei abläuft. Artefakte können, wie in Abbildung A.1 im Anhang dargestellt, bei Verwendung dieser Methode vor allem bei scharfen, farbigen Kanten entstehen, wenn diese durch einen gesubsampten Block verlaufen.

2.2. Diskrete Kosinus Transformation

Wie bereits oben beschrieben neigt der menschliche Sehsinn dazu niedrige Ortsfrequenzen eher zu erkennen, als höhere. Eine Ortsfrequenz ist definiert als „Anzahl bestimmter periodischer Erscheinungen bezogen auf einen räumlichen Abstand“ [Atm]. Wir erkennen also gröbere Strukturen mit einer niedrigen Ortsfrequenz eher als feinere Strukturen mit einer höheren. Um diesen Umstand nun auszunutzen muss das Ausgangsbild von der räumlichen Ebene auf eine Frequenzebene transformiert werden, damit anschließend, in dem darauf folgenden Schritt der Quantisierung, die höheren Frequenzen reduziert werden können. Diese Transformation lässt sich mittels einer zweidimensionalen Diskreten Kosinus Transformation (DCT) bewerkstelligen.

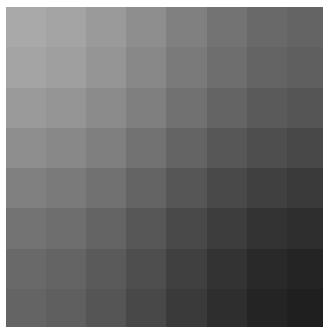
Die DCT ist eine Sonderform der Fouriertransformation, bei der eine Funktion mittels Sinusschwingungen approximiert wird. Die Fouriertransformation hat allerdings unter anderem den Nachteil, dass für jeden abgetasteten Punkt ein Tupel aus Amplitude und Phase bzw. Sinus und Kosinus Koeffizienten gespeichert werden muss. Die DCT nutzt nun den Umstand aus, dass das betrachtete Intervall begrenzt ist. Durch eine vertikale Spiegelung dieses Intervalls lassen sich die Sinus Anteile herauskürzen, wobei am Ende lediglich Kosinus Anteile übrig bleiben, also nur ein Koeffizient pro abgetasteten Punkt gespeichert werden muss. Des Weiteren bewirkt die Spiegelung, dass Start- und Endpunkt equivalent sind. Da die Fouriertransformation von einer unendlichen Folge ausgeht, muss der letzte Koeffizient den ggf. großen Unterschied zwischen Start- und Endpunkt ausgleichen. Sind diese Punkte aber equivalent, wird die Kraft des letzten Koeffizienten nicht verschwendet [Sym04]. Verarbeitet werden mit der zweidimensionalen DCT immer 8x8 Blöcke eines jeden Kanals mit der Formel:

2. Irrelevanzreduktion

$$F(u, v) = \frac{1}{4} C_u C_v \sum_{x=0}^7 \sum_{y=0}^7 f(x, y) \cos\left(\frac{(2x+1)u\pi}{16}\right) \cos\left(\frac{(2y+1)v\pi}{16}\right)$$

$$\text{wobei } \begin{cases} C_u = \frac{1}{\sqrt{2}} \text{ für } u=0, \text{ ansonsten } C_u=1 \\ C_v = \frac{1}{\sqrt{2}} \text{ für } v=0, \text{ ansonsten } C_v=1 \end{cases}$$

Die Abbildung 2.1 zeigt das Resultat einer angewandten DCT auf einen schwarz-weißen 8x8 Pixelblock, welcher aus jeweils einer horizontalen und einer vertikalen Kosinus Schwingung besteht. Der sogenannte DC Wert ist der erste Wert der Matrix und gibt die mittlere Helligkeit an. Alle anderen Komponenten beschreiben die relative Abweichung zu diesem Wert und werden gemeinhin als AC Werte betitelt, wobei diese zugleich die zum unteren rechten Rand hin höher werdenden Ortsfrequenzen repräsentieren. Wie bereits zu erkennen führt die DCT oftmals selbst schon durch Rundung auf ganzzahlige Ergebnisse zu einer Matrix mit einer erhöhten Anzahl gleicher Werte, die sich für die Anwendung weiterer, verlustfreier, Kompressionsmethoden eignet.



800	200	0	0	0	0	0	0
200	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0

Abbildung 2.1.: Mittels DCT gut komprimierbarer 8x8 Pixelblock
Links: Ausgangsbild, Rechts: Resultierende DCT-Matrix

2.3. Quantisierung

Im vorigen Schritt wurde durch Anwendung der Diskreten Kosinus Transformation eine Matrix mit den korrespondierenden Ortsfrequenzen eines 8x8 Pixelblocks gewonnen. Um nun tatsächlich eine Reduktion der höheren Ortsfrequenzen zu erreichen wird die Methode der Quantisierung angewandt. Hierbei wird eine ganzzahlige Division eines jeden DCT Koeffizienten mit einem Quantisierungswert vorgenommen. Das gerundete Ergebnis ist dann der quantisierte Wert. Durch diese Division und Rundung wird versucht die

bisher noch hohen Koeffizienten zu verkleinern, sowie in den höheren Frequenzbereichen möglichst auf Ergebnisse gleich Null zu kommen.

Im Fall von MPEG-1 wird hierfür ein Uniform Scalar Quantizer verwendet, bei dem die Eingangswerte durch Division der Schrittgröße auf Bereiche gleicher Größe abgebildet werden, wobei eine stufenähnliche Charakteristik entsteht. [Sym04]. Um die errechneten Ortsfrequenzen in Relation zur Wahrnehmung des menschlichen Auges zu reduzieren wird hierfür eine Quantisierungsmatrix verwendet. Diese beinhaltet separate Werte für jeden DCT Koeffizienten. Die Schrittgröße setzt sich für AC-Werte zusammen aus dem korrespondierenden Quantisierungswert der Quantisierungsmatrix und einem Quantisierungsfaktor (MQuant). Der Quantisierungsfaktor dient der Steuerung der Bildqualität und kann einen Wert zwischen 1 und 31 annehmen, wobei ein Quantisierungsfaktor von 1 für eine hohe Bildqualität sorgt, ein Faktor von 31 hingegen für eine stark reduzierte. Da das Auge sensibel gegenüber großräumigen Luminanzfehlern ist, wird der DC durch eine feste Schrittgröße von 8 dividiert. [11193] Eine Implementierung des vorgestellten Algorithmus ist in Listing A.1 im Anhang zu sehen.

In Abbildung A.2 des Anhangs ist die angewandte Quantisierung exemplarisch an einem Beispielbild mit der im MPEG-1 Standard voreingestellten Quantisierungsmatrix (siehe Anhang, Tabelle A.1) sowie Quantisierungsfaktoren von eins, 16 und 31 dargestellt. Bei höheren Quantisierungsfaktoren sind hier deutliche Qualitätsverluste zu erkennen, wobei die groben Strukturen des Bildes aber erhalten bleiben.

Durch die Anwendung der DCT wird ein eingehender 8x8 Pixelblock also in eine Darstellung transformiert, die es erlaubt mittels der Quantisierung vor allem enthaltene höhere Ortsfrequenzen zu reduzieren. Diese Prozesse führen zunächst jedoch nicht direkt zu einer Datenreduktion, da trotz des erhöhten Anteils gleicher Werte in der Matrix eben diese Werte auch gespeichert werden müssen. Allerdings wurde erreicht, dass die Entropiekodierung, welche im nachfolgenden Kapitel erläutert wird, bessere Kompressionsergebnisse erzielen kann.

3. Redundanzreduktion

Die Redundanzreduktion ist cool und sie reduziert Redundanz.

3.1. Entropiecodierung

3.2. Motion Compensation

Alle bis jetzt vorgestellten Ansätze der Videokompression beschäftigen sich mit der Kompression von Einzelbildern innerhalb eines Videos. Bei der Motion Compensation hingegen wird das Kompressionspotential ausgenutzt, dass innerhalb der Abhängigkeiten der Einzelbilder in einem Video steckt. Videos bestehen meist aus zusammenhängenden Szenen mit größtenteils unverändertem Inhalt innerhalb einer jeweils solchen Szene.

Man stelle sich zum Beispiel die folgende Szene vor: Eine statische Kamera filmt einen Mensch, unseren Protagonisten, der eine Straße entlang läuft und schließlich eine Bar betritt, eine typische Szene in Serien heutzutage.

Teilt man diese Szene in ihre Einzelbilder auf, stellt man schnell fest, dass die Einzige Bewegung der laufende Protagonist ist und der Hintergrund dabei komplett statisch verbleibt. Motion Compensation nutzt die Redundanz dieser statischen Hintergründe aus indem es diese jeweils nur ein Mal speichert und in den folgenden Bildern darauf referenziert um ein für den Zuschauer unverändertes Bild anzuzeigen. Da Videos üblicherweise zu großen Teilen mit statischen Bildteilen übersät sind, macht die von Motion Compensation erzielbare Kompression einen großen Teil des gesamten möglichen Kompressionspotentials innerhalb von Videos aus.

Damit Motion Compensation überhaupt funktionieren kann ist eine Aufteilung und Auswertung aller Video Einzelbilder (Frames) nötig.

3.2.1. Frames

Mit dem Kodieren teilt man alle Frames in eine bestimmte Bildart ein: Es gibt rein intra-codierte Frames, die sogenannten intracoded Frames (kurz I-Frames), bei denen es sich um einzelne Vollbilder, die Allein stehen und somit von keinem anderen Bild des Videos abhängen. Bei ihnen handelt es sich im Endeffekt einfach um ein für sich stehendes JPEG, was mit den üblichen Methoden der Bildkompression verkleinert wurde. Außerdem gibt es intercodierte Frames, die nur einen vorhergesagte Differenz des Inhalt in Abhängigkeit zu einem vorherigen I-Frame haben, die sogenannten predictive Frames (kurz P-Frames). Als letztes gibt es bipredictive Frames (kurz B-Frames), die sehr ähnlich zu P-Frames, die in zwei Richtungen intercodiert sind, nämlich indem sie die vorhergesagte Differenz des Inhaltes zum vorherigen I- oder P-Frame speichern. Um die Vorhersagung zu erreichen zu können wird eine Reihenfolge der Codierung gewählt, die ungleich der Reihenfolge der Anzeigereihenfolge ist, wie auf der Abbildung X erkennbar ist. Dadurch wird der sowieso schon komplexe Prozess zusätzlich erschwert.

Ein kompletter Szenenwechsel, also das Ändern des kompletten Bildes, ohne statische Zusammenhänge, muss dem Encoder immer mitgeteilt werden. Dieser muss dann einen neuen I-Frame codieren, auf dem die folgenden P- und B-Frames basieren. Dadurch wird die potentielle Gefahr einer starken Artefaktbildung vorgebeugt.

Wenn man diese Aufteilung jedoch jeweils nur einmal pro Szene anwenden würde, würden mehrere Probleme bei wahllosem Zugriff entstehen. Wenn der I-Frame einer Szene fehlt oder übersprungen wird, würden die Änderungen, die in den folgenden P- und B-Frames festgehalten wurden, auf den falschen I-Frame angewendet, sodass im Video starke Artefakte entstehen. Beim Ausfall eines P-Frames einer Szene würde grundsätzlich das Gleiche geschehen, jedoch nur bei den noch folgenden P-Frames der Szene.

Um diese unschönen Artefakte beim Vor- und Zurückspulen zu verhindern, dürfte nur zu einem I-Frame gesprungen werden, welches bei einer Aufteilung pro Szene jeweils der Anfang einer neuen Szene wäre.

Da bei einem Großteil der Anwendungsfälle von Videos jedoch eine fast vollständig wahlfreier Zugriff gewünscht ist, teilt man sie in viele kleine aufeinanderfallende Bildergruppen (Group of pictures, kurz GOP) auf. Eine GOP wird meist mit 2 Parametern angegeben, zum Beispiel N und M. Dabei ist N eine bestimmte Anzahl von Frames aus denen die GOP besteht, also die Distanz von einem I-Frame zum nächsten I-Frame. M gibt die Distanz von einem I- oder P-Frame, bis zum jeweils Nächsten an, somit ist M-1 die Anzahl von

3. Redundanzreduktion

B-Frames, die nach einem I- oder P-Frame folgen. Eine Bildergruppe fängt immer mit einem I-Frame an und wiederholt sich bis zum Ende eines Videos mit einem konstanten Schema.

Mit den Parametern $N=12$ und $M=4$, würde die GOP dann aussehen wie auf der Abbildung X. (IBBBPBBBBPBBB I...) TODO: ABBILDUNG

Bei MPEG ist eine Aufteilung mit den Parametern $M=3$ bis 4 und $N= 11$ bis 15 üblich.

Betrachtet man eine übliche Framerate von 25 ist somit wahlfreier Zugriff bis auf die Hälfte einer Sekunde gegeben. Außerdem wird dadurch bei leichten Übertragungsfehlern einer Videodatei der Schaden minimiert, sodass das vom Endnutzer gesehene Bild nur maximal eine halbe Sekunde Artefakte anzeigt.

3.2.2. Makroblocks

TODO: Jeder Interodierte Frame wird in sogenannte Makroblöcke unterteilt. Diese Makroblöcke werden beim Codieren zum Vergleichen mit dem vorher codierten Bild mittels Block matching algorithmus benutzt.

3.2.3. Motion Estimation and Compensation

: TODO: Wenn ein ähnlicher Block gefunden wird, wird der Block mittels dem resultierenden Motion Vektor encodiert.

4. Ausblick

ÄÖÜäöüß

5. Zusammenfassung

ÄÖÜäöüß

A. Weitere Abbildungen und Tabellen

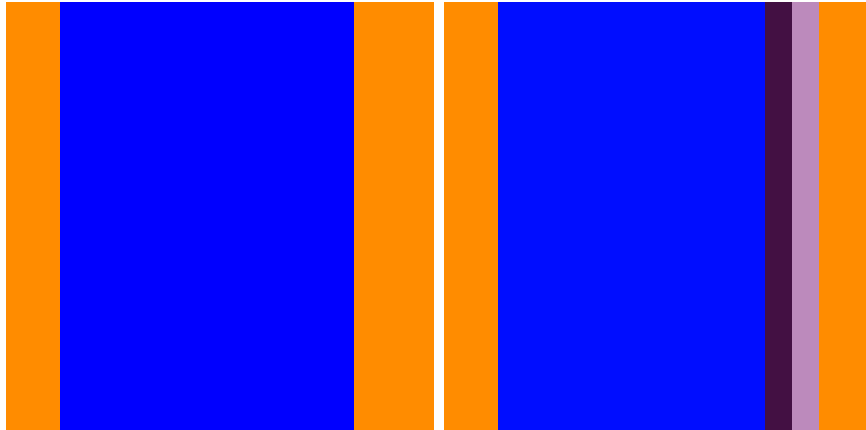


Abbildung A.1.: Artefakte durch Chroma Subsampling

Links: Original, Rechts: Subsampled. Die rechte Kante des blauen Farbblocks liegt in gesubsampten 2x2 Blöcken, wodurch Artefakte entstehen. Die linke Kante liegt zwischen zwei 2x2 Blöcken, weshalb es zu keiner falschen Darstellung kommt.

```
def quantize(dct, quantizer, MQuant=1):
    result = np.empty_like(dct)
    for x, row in enumerate(dct):
        for y, coefficient in enumerate(row):
            if x == 0 and y == 0:
                result[x][y] = int(coefficient / 8)
            else:
                result[x][y] = int( 8 * coefficient / (MQuant *
                    quantizer[x][y] ))
    return result
```

Listing A.1: Implementierung des Quantisierungsprozesses nach MPEG-1 Standard ohne Clipping

A. Weitere Abbildungen und Tabellen



Abbildung A.2.: Ergebnis der Quantisierung mit verschiedenen Quantisierungsfaktoren
Oben links: Original, Oben rechts: Quantisiert mit Faktor 1, Unten links: Quantisiert mit Faktor 16, Unten rechts: Quantisiert mit Faktor 31.

Mit zunehmendem Quantisierungsfaktor ist ein ansteigender Verlust der Bildqualität zu beobachten, wobei grobe Strukturen weitestgehend erhalten bleiben. Original nach [Cag16]

8	16	19	22	26	27	29	34
16	16	22	24	27	29	34	37
19	22	26	27	29	34	34	38
22	22	26	27	29	34	37	40
22	26	27	29	32	35	40	48
26	27	29	32	35	40	48	58
26	27	29	34	38	46	56	69
27	29	35	38	46	56	69	83

Tabelle A.1.: Voreingestellte MPEG-1 Intracoding Quantisierungsmatrix. [Sym04]

RLE	Genutzte Optionen		Größe in Kilobyte	Ratio
	Chroma Subsampling	Quantisierung mit Faktor		
X		-	258.38	3.76
X	X	-	145.95	6.66
X	X	1	58.32	16.67
X	X	16	18.59	52.29
X	X	31	16.23	59.89

Tabelle A.2.: Testergebnisse der angewandten Kompressionsalgorithmen bei einer Ausgangsgröße von 970 Kilobyte des Originalbildes

Literatur

- [11193] ISO/IEC 11172-2:1993. *Information technology – Coding of moving pictures and associated audio for digital storage media at up to about 1,5 Mbit/s – Part 2: Video*. Standard. International Organization for Standardization, 1993.
- [Akr14] Shahriar Akramullah. *Digital Video Concepts, Methods, and Metrics*. Berkeley, CA: Apress, 2014, S. 19.
- [Atm] AtmWiki. *Ortsfrequenz*. URL: <http://www.otterstedt.de/wiki/index.php/Ortsfrequenz> (besucht am 02.12.2016).
- [Cag16] Brooke Cagle. 2016. URL: <https://unsplash.com/photos/EBhyjAclIPo> (besucht am 03.12.2016).
- [Int95] International Telecommunication Union, Telecommunication Standardization Sector (ITU-T). *Recommendation ITU-R BT.601-5: Studio encoding Parameters of digital television for standard 4:3 and wide-screen 16:9 aspect ratios*. 1995.
- [Poy] Charles Poynton. *Chroma subsampling notation*. URL: http://www.poynton.com/PDFs/Chroma_subsampling_notation.pdf (besucht am 30.11.2016).
- [Sym04] Peter Symes. *Digital Video Compression*. New York: McGraw-Hill, 2004, S. 76–79, 94–95, 162.