

Prof. Dr. Stefan Weinzierl

Musterlösung: 4. Aufgabenblatt

1. Aufgabe: Dynamik und SNR der Zahlenformate

Generieren Sie eine Periode eines Sinussignals mit 1024 samples/Periode und quantisieren Sie es auf eine Wortbreite von 8 bit im Festkomma-Format, sowie auf 16 bit im Fließkomma-Format. Hierbei betrage nach IEEE 754r die Wortbreite der Mantisse 10 bit und die des Exponenten 5 bit.

Gleitkomma-Quantisierung¹

Eine Gleitkomma-Zahl wird dargestellt in der Form:

$$x_Q = M \cdot 2^E$$

Die Mantisse M wird dabei zur eindeutigen Darstellung normalisiert und nimmt nur Werte zwischen 0,5 und 1 (ohne 1!) an. Sie wird im Festkomma-Format dargestellt mit einer Wortbreite von w_m bit für die Zahl selbst sowie 1 bit für das Vorzeichen.

Für den Exponenten gilt $E = e - bias$. e ist eine Zahl zwischen 1 und $2^{w-e} - 2$ (hier werden die Sonderfälle 0 und Infinity/NaN berücksichtigt). Um auch negative Exponenten zu ermöglichen wird der so genannte $bias$ abgezogen, der sich wie folgt berechnen lässt (und auch in der IEEE-Norm definiert ist): $bias = e_{max}/2$.

Die Werte für 0 und unendlich sind als Sonderfälle definiert (siehe KT 2-Skript), hierbei ist zu beachten, dass sich ein positiver und ein negativer Wert für 0 ergeben.

```
% Eingangssignal mit N samples
N = 1024;
n = 0:1/N:1-1/N;
y = sin(2*pi*n);

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%% fixed point %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

% ---- fixed point Quantisierung- --- %
w = 8;
q = 2/2^w;
y_fixed = quant(y,q);

% die hoechste Quantisierungsstufe (1) wird weggelassen, um genau 2^w Stufen zu
% erreichen
maximum = max(y);
for i = 1:length(y_fixed);
    if (y_fixed(i) > (maximum - q));
        y_fixed(i) = (maximum - q);
    end
end

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%% floating point %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

s = [1 0]; % Vorzeichenbit
w_e = 5; % Wortbreite des Exponenten
```

¹ siehe auch: Zölzer (2005), Abschnitt „Zahlendarstellung“

```

w_m = 10; % Wortbreite der Mantisse
e = 1 : 2^w_e - 2; % Vektor mit allen Werten von e
bias = max(e)/2; % bias, wird von e abgezogen
q_m = 2^-w_m; % Quantisierungsstufe der Mantisse (diese wird im fixed point
Format dargestellt)
m = 0.5:q_m:1-q_m; % Vektor mit allen Werten der Mantisse

% Vorbereitung für die for-Schleife
l_s = length(s);
l_e = length(e);
l_m = length(m);

% Es wird eine Vektor mit allen Werten des floating point Formats erzeugt
for i = 0 : l_s - 1;
    for j = 0 : l_e - 1;
        floatingpoint_(1 + j*l_m : l_m + j*l_m) = (-1)^s(i+1) * m .* 2 ^ (e(j+1) -
bias);
    end
    floatingpoint(1 + i*(l_e*l_m) : l_e*l_m + i*(l_e*l_m)) = floatingpoint_;
end

floatingpoint = [floatingpoint 0 -0]; % es wird der Sonderfall 0 eingeführt
floatingpoint = sort(floatingpoint); % die Werte werden nach ihrer Groesse sortiert

% --- floating point Quantisierung ---- %
% Das Eingangssignal wird auf das Maximum des floating point Vektors
% skaliert
y_floating = y*max(floatingpoint);

% Um zu ermitteln, welcher Wert des floating point Vektors den
% Eingangssignalwerten am nächsten ist (also auf welche floating point Zahl
% die Eingangswerte gerundet werden), wird die Differenz zwischen
% floating point Zahlen und Eingangssignalwerten gebildet (differenz). Die
% Position der kleinsten Differenz wird ermittelt (index_min), um den
% Eingangssignalwert auf eben diese Stufe im floating point Vektor zu runden.
for i = 1:length(y);
    differenz = floatingpoint - y_floating(i);
    [min_differenz index_min] = min(abs(differenz));
    y_floating(i) = floatingpoint(index_min);
end

% Das quantisierte Signal wird wieder auf die ursprüngliche Skalierung
% gebracht.
y_floating = y_floating/max(floatingpoint);

```

a) Wie groß ist die Dynamik für beide Signale in dB?

Die Dynamik lässt sich berechnen, in dem man das Verhältnis zwischen maximal möglicher Amplitude und minimal möglicher Amplitude des Signals bildet:

$$DR = 20 \log \left(\frac{x_{\max}}{x_{\min}} \right)$$

Im Falle der Festkomma-Quantisierung (mit midtread-Kennlinie) ist ein Zahlenbereich von -2^7 bis $2^7 - 1$, bzw. $[-128 \dots 127]$, darstellbar. Dieser wird durch die Vollausslenkung des Sinus ausgenutzt. Die minimale Auslenkung liegt zwischen ± 1 .

Es ist oberste Quantisierungsstufe

$$x_{\max} = 1 - 2^{1-w} = 1 - q$$

und die kleinste Quantisierungsstufe

$$x_{\min} = 2^{1-w} = q.$$

Dies resultiert in einer Dynamik von

$$DR = 20 \log \left(\frac{1-q}{q} \right) = 20 \log \left(\frac{1-2^{-7}}{2^{-7}} \right) = 20 \log \left(\frac{2^7-1}{1} \right) = 20 \log \left(\frac{127}{1} \right) = 42,08 \text{ dB}$$

Im Falle der Fließkomma-Quantisierung ist

$$x_{\max} = M_{\max} \cdot 2^{E_{\max}} \text{ und } x_{\min} = M_{\min} \cdot 2^{E_{\min}},$$

dabei ist die maximale Mantisse

$$M_{\max} = (1 - 2^{-w_m}) = 1 - q_m.$$

Die minimale Mantisse M_{\min} ist wegen der Normalisierung gleich 0,5.

Der Zahlenbereich des Exponenten liegt zwischen [--14...15]:

$$E_{\min} = 1 - bias = -2^{w_e - 1} + 2 \text{ und } E_{\max} = (2^{w_e} - 2) - bias = 2^{w_e - 1} - 1.$$

Dies führt auf eine maximale Auslenkung zwischen $\pm 32'736$. Die minimale Auslenkung liegt zwischen $\pm 3,052 \cdot 10^{-5}$.

$$DR = 20 \log \left(\frac{M_{\max} \cdot 2^{E_{\max}}}{M_{\min} \cdot 2^{E_{\min}}} \right) = 20 \log \left(\frac{(1 - q_m) \cdot 2^{2^{w_e} - bias}}{0,5 \cdot 2^{1 - bias}} \right) = 20 \log \left(\frac{(1 - 2^{-10}) \cdot 2^{15}}{0,5 \cdot 2^{-14}} \right) = 180,61 \text{ dB}$$

$$DR_{\text{fixed}} = 20 \cdot \log_{10}((1-q)/q)$$

$$DR_{\text{floating}} = 20 \cdot \log_{10}((\max(m) \cdot 2^{(\max(e)-bias)}) / (\min(m) \cdot 2^{(\min(e)-bias)}))$$

b) Transformieren Sie beide Signale mittels FFT in den Frequenzbereich und berechnen Sie jeweils den SNR. Wie interpretieren Sie die Ergebnisse im Hinblick auf die in b) berechnete Dynamik?

Der SNR ist das Verhältnis zwischen Leistung des Nutzsignals und Leistung des Fehlers:

$$SNR = 10 \log \left(\frac{P_S}{P_N} \right)$$

Eine einfache Berechnung des SNR basierend auf der FFT-Technik ist dann möglich, wenn die Abtastrate ein Vielfaches der Sinusperiode ist bzw. Signallängen verwendet werden, die mit einer ganzzahligen Anzahl von Perioden in das FFT-Analyseintervall passen. Genau diese Voraussetzung wurde mit den Rahmenbedingungen in dieser Aufgabe geschaffen (das Signal passt mit genau einer Periode in einen Block der FFT). In diesem Fall sind alle Rauschanteile in den diskreten harmonischen des FFT-Spektrums enthalten, sodass zur Berechnung des SNR folgende Formel verwendet werden kann:

$$SNR = 10 \log \left(\frac{y_1^2}{y_2^2 + y_3^2 + y_4^2 + \dots + y_{N/2-1}^2} \right), \text{ mit } y = \text{Amplituden der Spektralkomponenten}$$

Hierbei wird die Parsevalsche Beziehung ausgenutzt:

```
% FFT der quantisierten Signale
Y_fixed = fft(y_fixed);
Y_floating = fft(y_floating);

% Betragsbildung und Beschränkung auf N/2 samples
Y_fixed = abs(Y_fixed(1:N/2));
Y_floating = abs(Y_floating(1:N/2));

%----- SNR-Berechnung -----%

% fixed point
y_grund_fixed = Y_fixed(2); % Grundschiwingung
P_teil_fixed = sum(Y_fixed(3(end-1)).^2); % Gesamtleistung der Oberschwingungen
```

```

P_grund_fixed = y_grund_fixed^2;           % ohne DC und Grundschiwingung
                                           % Leistung der Grundschiwingung

% Der SNR ist das Verhaeltnis der Leistung der Grundschiwingung zur Leistung
% der Oberschiwingungen.
snr_fixed = 10*log10(P_grund_fixed/P_teil_fixed)

% floating point
y_grund_floating = Y_floating(2);          % Grundschiwingung
P_teil_floating = sum(Y_floating(3:(end-1)).^2); % Gesamtleistung der Oberschiwin-
gungen ohne DC und Grundschiwingung
P_grund_floating = y_grund_floating^2;     % Leistung der Grundschiwingung

% Der SNR ist das Verhaeltnis der Leistung der Grundschiwingung zur Leistung
% der Oberschiwingungen.
snr_floating = 10*log10(P_grund_floating/P_teil_floating)

```

Die Berechnung liefert einen SNR von 49,29 dB für die 8-Bit-Festkommazahl und 69,72 dB für die 16-Bit-Fließkommazahl. Dies zeigt, dass obwohl im Fließkommaformat eine deutlich größere Dynamik erreicht werden kann, der SNR bei einer Festkomma-Kodierung mit derselben Wortbreite sogar kleiner wäre. Dies liegt daran, dass die Quantisierungsintervalle bei großen Amplituden ebenfalls immer größer werden, sodass natürlich auch der relative Fehler größer wird.

c) Welche Frequenz müsste das Signal bei einer Abtastfrequenz von 48 kHz haben?

$$\text{Signalfrequenz} = \frac{\text{Abtastfrequenz}}{\text{samples in einer Periode}} = \frac{48000 \text{ Hz}}{1024} = 46,875 \text{ Hz}$$

Diese Frequenz müsste für das Testsignal gewählt werden, um bei einer FFT-basierten Messung des SNR eines A/D-Wandlers ein exaktes Ergebnis zu erhalten.

2. Aufgabe: A/D- und D/A-Wandler

Die Datenblätter zweier digitaler Aufzeichnungsgeräte mit 24 bit Wortbreite weisen folgende Werte auf:

Datenblatt 1:

THD+N (A/D): <.001% typ @ 1 kHz, @ clip level -0.5 dB

Dynamic Range (A/D): 109 dBA

THD+N (D/A): <.003% typ @ 1 kHz, @ clip level -0.5 dB

Dynamic Range (D/A): 111 dBA

Datenblatt 2:

Dynamic range: 103dB A/D, >103dB D/A (A-weighted). THD+N: -91 dB.

a) Erläutern Sie das Messverfahren für die Größen THD+N und Dynamic Range. Welche Details über die verwendeten Messsignale können Sie den obigen Angaben entnehmen?

THD+N :

THD+N (total harmonic distortion plus noise) ist das Verhältnis aller Oberschwingungen *inkl.* Grundrauschen zum voll ausgesteuerten Messsignal (Effektivwerte).² Der THD+N ist also einfacher zu messen als der THD allein, da bei letzterem auch Aliasinganteile der Obertöne an völlig anderen Positionen als der Obertonreihe auftreten können und berücksichtigt werden müssten. Diese aufwändige Spektralanalyse kann entfallen, da beim THD+N einfach „alles“ im hörbaren Audiobereich in den Messwert einfließt, was sich vom Messton unterscheidet. Die Messung berücksichtigt somit nicht nur harmonische Oberwellen, sondern das gesamte Störspektrum einschließlich unharmonischer Anteile, Einstreuungen, Brummen, Rauschanteile u.ä.

Ablauf der Messung:

- rege mit $f_0 = 997\text{Hz}$ Ton bei $-0,5\text{ dB FS} / -1\text{ dB FS}$ an
- messe Effektivwert von Rauschen und Verzerrungsprodukten im gesamten Audioband bei Unterdrückung von f_0 durch Notch-Filter
- bilde Verhältnis zur FS-Spannung und berechne Pegel

Notation z.B.: THD+N (997 Hz, -1 dB FS) = -85 dB FS

Für das Grundrauschen wird üblicherweise kein Gewichtungsfiler verwendet, da von voll ausgesteuerten also „lauten“ Nutzsignalen ausgegangen wird (ungefähr linearer Bereich der Hörkurve).

Dynamic Range:

Dynamic Range bezeichnet das Verhältnis des (frequenzbewerteten) Grundrauschens in Anwesenheit eines (leisen) Signals zur Vollaussteuerung des Wandlers (in [dB FS]). Es erscheint in der AES-17-Norm als „*Signal-to-Noise Ratio (SNR) in the presence of a signal*“.

Ablauf der Messung:

- zur Angabe des DR-Wertes muss vorher eine FS-Messung stattgefunden haben
- dann wird ein Messton $f_0 = 997\text{ Hz}$ bei -60 dB FS benutzt, damit möglichst keine Nichtlinearitäten angeregt werden, der Wandler aber sicher aktiv ist (keine Stummschaltung)
- f_0 ausfiltern, THD+N Messung relativ zur FS-Spannung durchführen

Notation z.B.: DR = 85 dB FS A

- Gewichtungsfiler z.B.: A, CCIR
- FS bezieht sich auf den Output FS-Wert

Detailinformationen über vorliegende Testsignale:

1. Wandler: 1kHz, $-0,5\text{ dB FS}$
2. Wandler: ??? k.A.

² Die Festlegung der Fullscale-Amplitude erfolgt dabei während der Messung selbst, da diese nach der AES-17 Norm wie folgt definiert ist: „... wenn das digitale Signal nicht zugänglich ist, wird Input-Fullscale $0,5\text{ dB}$ unterhalb des Amplitudenwertes eines 997-Hz -Sinustons gesetzt, bei dem 1% THD+N oder $0,3\text{dB}$ Kompression erreicht werden...“.

Die AES-17 Norm schreibt eigentlich die Benutzung eines 997-Hz-Tons vor, da dieser zu viel mehr unterschiedlichen Codewörtern führt als ein 1000-Hz-Ton (bei 48 kHz Abtastrate z.B. nur 48), und sonst je nach Phasenlage zu Beginn der Abtastung evtl. nie ein Abtastwert beim Amplitudenmaximum erreicht wird!

b) Wie weit weichen die gemessenen Werte von denen eines idealen 24-bit Wandlers ab?

Der theoretisch erreichbare SNR für Sinussignale auf Grund von Quantisierungsrauschen lautet:

$$\text{SNR}_{\text{Sinus}} = n \cdot 6,02 \text{ dB} + 1,76 \text{ dB} = 146,2 \text{ dB}$$

THD+N und DR müssten dem theoretischen SNR entsprechen.

Festzustellende Abweichungen:

DR:

Bewertete Messungen können mit der theoretischen Berechnung (unbewertet) nicht direkt verglichen werden (die A-Bewertung erzeugt bei Messungen an Audioprozessoren im Schnitt 2 dB bessere SNR-Werte).

THD+N:

Erster Wandler : Eingang: THD+N < .001% $\rightarrow 20 \cdot \log_{10}(0.00001) = -100 \text{ dB FS}$
 $\rightarrow 46 \text{ dB schlechter}$
Ausgang: THD+N < .003% $\rightarrow -90,5 \text{ dB FS}$
 $\rightarrow 55,7 \text{ dB schlechter}$

Zweiter Wandler : -91 dB $\rightarrow 55,2 \text{ dB schlechter}$

Hier ist das Messsignal unbekannt, die theoretisch erreichbaren 146,2 dB beziehen sich aber auf reine Sinustöne.

c) Worauf sind diese Abweichungen zurückzuführen? Geben Sie einige potentielle Fehlerquellen digitaler Übertragungssysteme an. Wieso unterscheiden sich die Werte für THD+N und DR?

- bei 24 bit Auflösung bestimmen vor allem thermale Rauschquellen (z.B. der analogen Eingangskomponenten, Übertrager) die Rauschleistung des Wandlers
- induzierte Störungen aufgrund mangelhaften Erdungslayouts auf der Wandlerplatine
- herstellerseitige Dynamikbeschränkung für Sicherheits-Headroom (z.B: 0 dB der Anzeige sind tatsächlich -12 dB FS)
- Begrenzung durch analoges Vorverstärkerequipment (Sättigung), Limiter, Softclipper
- Schwankungen der Referenzspannung oder der Spannungsversorgung
- Amplitudenfehler: Verschiebung des Nulldurchgangs, Verstärkungsfehler, differentielle/integrale Nichtlinearitäten, Glitches
- Zeitfehler: Jitter im A/D-D/A-Kreis
- Mängel der Antialiasing-/ Rekonstruktionsfilter

Gründe für Unterschiede bei DR / THD+N:

- bei Systemen mit wenig Verzerrungsprodukten sind THD+N und DR annähernd gleich, es messen dann beide hauptsächlich den Rauschpegel
- THD+N ist aber meist schlechter, da sich durch den hohen Messsignalpegel (z.B. -0.5 dB FS) deutliche Verzerrungsprodukte ergeben
- Frequenzbewertung bei der DR-Messung

d) Vergleichen Sie anhand der Daten die Qualität der verwendeten Wandler.

Erster Wandler:

A/D-Strecke:

- THD+N = -100 dB
- DR = 109 dB, ist 9 dB besser als THD+N

D/A-Strecke:

- THD+N = -90,5 dB
- DR = 111dB, ist 20.5 dB besser als THD+N. 20 dB Differenz bedeuten, dass THD+N bei dieser Messung von Verzerrungsprodukten dominiert wurde.

→ Messsignalpegel ist mit -0.5 dB FS sehr hoch

→ THD+N Differenz zwischen A/D- und D/A-Strecke 9,5 dB

Zweiter Wandler:

A/D-Strecke:

- THD+N = -91 dB; A/D-Strecke deutlich schlechter als erster Wandler
- DR = 103 dB A; 6 bzw. 8 dB niedriger als erster Wandler

Da hier weder Signalpegel noch -frequenz gegeben sind, sind die beiden Wandler aber schlecht vergleichbar.