



Akademia Górniczo-Hutnicza im. Stanisława Staszica w Krakowie

Wydział Informatyki

PRACA DYPLOMOWA

Generacja muzyki przy pomocy dużych modeli językowych

Music generation with Large Language Models

Autor:	Filip Ręka
Kierunek:	Informatyka — Data Science
Opiekun pracy:	dr hab. Maciej Smółka prof. AGH

Kraków, 2024

Tutaj możesz umieścić treść podziękowań. Tutaj możesz umieścić treść podziękowań. Tutaj możesz umieścić treść podziękowań. Tutaj możesz umieścić treść podziękowań.

Streszczenie

Duże modele językowe (ang. *Large Language Models* **LLM**) charakteryzują się zdolnością do generacji języka oraz innych zadań w przetwarzania języka naturalnego, takich jak na przykład klasyfikacja. Zdolność tą nabierają podczas czasochłonnego oraz intensywnego obliczeniowego treningu metodami samo odraz pół-nadzorowanego, podczas którego uczą się one relacji z wielkiej ilości dokumentów tekstowych. LLMy mogą zostać wykorzystane do generacji tekstu, formy generatywnej sztucznej inteligencji, poprzez pobieranie tekstu wejściowego i wielokrotne przewidywanie kolejnego tokenu lub słowa w tekście. Strukturę muzyki można porównać struktury tekstu pisanego, gdzie każda nuta odpowiada literze lub słowu, akordy zdaniom a dłuższe i sekwencje paragrafom. Poniższa praca, zamierza zbadać możliwości generacyjne LLMów wytrenowanych na muzycznych zbiorach danych.

Abstract

Abstract in English [\[1\]](#) ...

Spis treści

Lista kodów źródłowych	xiii
1 Wstęp	1
1.1 Motywacja	1
1.2 Cel pracy	1
1.3 Zarys pracy dyplomowej	2
2 Opis aktualnego stanu wiedzy	3
2.1 Podobieństwa pomiędzy muzyką a tekstem	3
2.2 Cyfrowa reprezentacja muzyki	3
2.2.1 WAV i MP3	4
2.2.2 MIDI (ang. <i>Musical Instrument Digital Interface</i>)	5
2.2.3 Notacja ABC	6
2.2.4 Porównanie zapisu muzycznego	9
2.3 Zbiory danych	11
2.3.1 Johann Sebastian Bach Chorales	11
2.3.2 The MAESTRO v3.0	11
2.3.3 IrishMAN	12
2.4 Architekturze transformera	13
2.4.1 Sieci RNN	13
2.4.2 LSTM	14
2.4.3 Algorytm uwagi (ang. <i>attention</i>)	15
2.4.4 Transformer	16
2.5 Architektura <i>state space</i>	19
2.5.1 Mamba	19
2.5.2 Mamba	22
3 Propozycja rozwiązania	25
4 Przeprowadzenie eksperymentów	27
4.1 Prezentacja otrzymanych wyników	27

4.2	Sposoby oceny wyników	27
4.3	Muzyczny test Turinga	29
4.3.1	Testowanie melodii innymi LLM-ami	30
5	Zakończenie	33
	Bibliografia	35
5		

Zawartość spisu treści — tytuły rozdziałów oraz ich liczba zależą od tematyki pracy — należy ustalić z opiekunem pracy.

Spis rysunków

2.1	Fragment chorału J.S. Bacha	4
2.2	Plik .wav otworzony w programie Audacity.	5
2.3	Muzyka zapisana w pliku MIDI	6
2.4	Zapis wielu głosów w notacji ABC.	7
2.5	Prosta sekwencja muzyczna	9
2.6	Reprezentacja muzyki w postaci tokenów	9
2.7	Przykład utworu pochodzącego ze zbioru MAESTRO.	12
2.8	Przykład wielogłosowego fragmentu ze zbioru IrishMAN.	13
2.9	Schemat budowy fragmentu sieci RNN.	13
2.10	Komórki LSTM połączone między sobą.	14
2.11	Schemat transformera.	17
2.12	Algorytm uwagi z maską.	18
2.13	Schemat modelu przestrzeni stanów.	20
2.14	Schemat modelu Mamba.	22
4.1	Wynik modelu.	28
4.2	Wygenerowana melodia.	28
4.3	Pytanie i odpowiedź modelu ChatMusician.	31
4.4	Odpowiedź modelu na fałszywy fragment.	32
4.5	Specjalnie stworzony “brzydki” fragment muzyki.	32

Spis tabel

2.1 Porównanie różnych zapisów muzyki	9
---	---

Lista algorytmów

Lista kodów źródłowych

1. Wstęp

Muzyka od zawsze stanowiła istotny element ludzkiego życia, inspirując, emocjonując i łącząc ludzi na różnych poziomach. Tradycyjnie proces tworzenia muzyki był zarezerwowany dla utalentowanych muzyków i kompozytorów, którzy posiadali wyjątkowy dar tworzenia dźwięków i melodii. Jednakże, w erze cyfrowej i rozwoju sztucznej inteligencji, pojawiają się nowe możliwości w dziedzinie generacji muzyki. Modelowanie generatywne, będące obszarem sztucznej inteligencji, pozwala na tworzenie nowych danych, w tym również muzyki, na podstawie wzorców i reguł wykrytych w zbiorze treningowym. Dynamiczny rozwój dziedziny przetwarzania języka naturalnego (*NLP*) ukazuje skuteczność tego podejścia w problemach generacji danych sekwencyjnych, do których należy właśnie tekst oraz muzyka.

1.1. Motywacja

Przez ostatnie lata widzimy szybki rozwój dużych modeli językowych (*LLM*), które w szczególności dobry sposób radzą sobie z rozumieniem tekstu w wielu językach. Modele te uczone są na wielkich korpusach danych, przez co są w stanie nauczyć się słownictwa oraz zasad gramatyki, a następnie na podstawie *promptu* udzielonego przez użytkownika wygenerować odpowiednią odpowiedź. Modele te znalazły zastosowanie w maszynowej translacji tekstu, analizie sentymentu, odpowiadaniu na pytania użytkownika czy nawet generacji kodu na podstawie poleceń oraz kontekstu nawet dużych projektów. Muzyka swoją strukturą jest bardzo podobna do tekstu. Każda nuta, tak jak słowo, nie tylko jest zależna od tych, które występują przed nią, ale również od szerszego kontekstu utworu jak również i tonacji w której dany utwór został napisany. Wykonując to porównanie, można łatwo zauważyć, dlaczego modele analizujące tekst oraz języki pisane są potencjalnie dobrymi kandydatami do próby analizy i generacji muzyki.

1.2. Cel pracy

Celem pracy jest przyjrzenie się dużym modelom językowym oraz ich zastosowaniu w celu generacji muzyki. Zawierać będzie opis modeli, które obecnie są wykorzystywane w architekturach *LLM*-ów oraz analizę odpowiadającą na pytanie dlaczego dane rozwiązanie sprawdza się w przedstawionym w pracy problemie. Następnie zostanie przedstawiony problem analizy

oraz generacji muzyki oraz jej możliwe cyfrowe reprezentacje. Zostaną również przedstawione dostępne zbiory danych, na podstawie których modele mogą zostać trenowane. Praca przedstawi wyniki użytych metod w celu ustalenia, który model i która reprezentacja muzyki najlepiej nadaje się do problemu przedstawionego w pracy.

1.3. Zarys pracy dyplomowej

Uwaga 1.1. Tutaj napiszę dopiero jak będzie więcej treści bo chyba troszkę bez sensu na razie pisać o strukturze jak jej jeszcze nie ma "_(ツ)_/"

2. Opis aktualnego stanu wiedzy

2.1. Podobieństwa pomiędzy muzyką a tekstem

Struktura tekstu i muzyki wykazuje wiele podobieństw, co pozwala na ich analizę za pomocą podobnych narzędzi i metod. Muzyka tak jak i tekst składa się z sekwencji znaków następujących w pewnej sekwencji. W obu przypadkach pojawiają się pewne schematy, czym w przypadku muzyki jest pewna progresja akordów lub powtarzający się fragment, a w tekście pewne związki frazeologiczne, powtarzające się wyrażenia lub ściśle określone reguły ortograficzne i składniowe. Te reguły w kontekście muzycznym są odpowiednikiem harmonii oraz strukturze i poprawności zapisu muzycznego. Sam sposób generacji sekwencji rozumianych przez komputer z muzyki zostanie poruszony w rozdziale 2.2.

Generacja muzyki w zasadzie nie różni się od generacji tekstu. Oba modele działają autoregresyjnie. Model autoregresyjny (AR) jest modelem, który modeluje pewne zjawisko sekwencyjne, a następnie iteracyjne na podstawie poprzednich sekwencji, przewiduje kolejny element, dokłada go do przewidzianej sekwencji, i ponownie “przepuszcza” ją przez model z dodatkowym elementem. Wyjściem modelu podczas jego pojedynczego kroku jest rozkład prawdopodobieństwa nad możliwymi tokenami zazwyczaj obliczanego z pomocą funkcji *softmax*.

$$\begin{aligned} \text{sekwencja wejściowa: } x_1, x_2, x_3, \text{ gdzie } x_i \in X \\ P_{t1}(X) = \text{Model}(x_1, x_2, x_3) \\ x_4 \sim P_{t1}(X) \\ P_{t2}(X) = \text{Model}(x_1, x_2, x_3, x_4) \\ x_5 \sim P_{t2}(X) \\ \dots \end{aligned} \tag{2.1}$$

Przykładowe podejście do autoregresyjnej generacji sekwencji pokazano w równaniu 2.1. Algorytm jest agnostyczny w stosunku do typu danych, więc w X mogą być zarówno wszystkie słowa, litery lub elementy notacji muzycznej.

2.2. Cyfrowa reprezentacja muzyki

W kulturze zachodniej nuty zapisuje się na pięcioliniach ułożonych jedna pod drugą. Jeśli muzyka jest wielogłosowa typowo pięciolinie są okalane przy pomocy nawiasów klamrowych.

Chorał no. 7

J.S. Bach



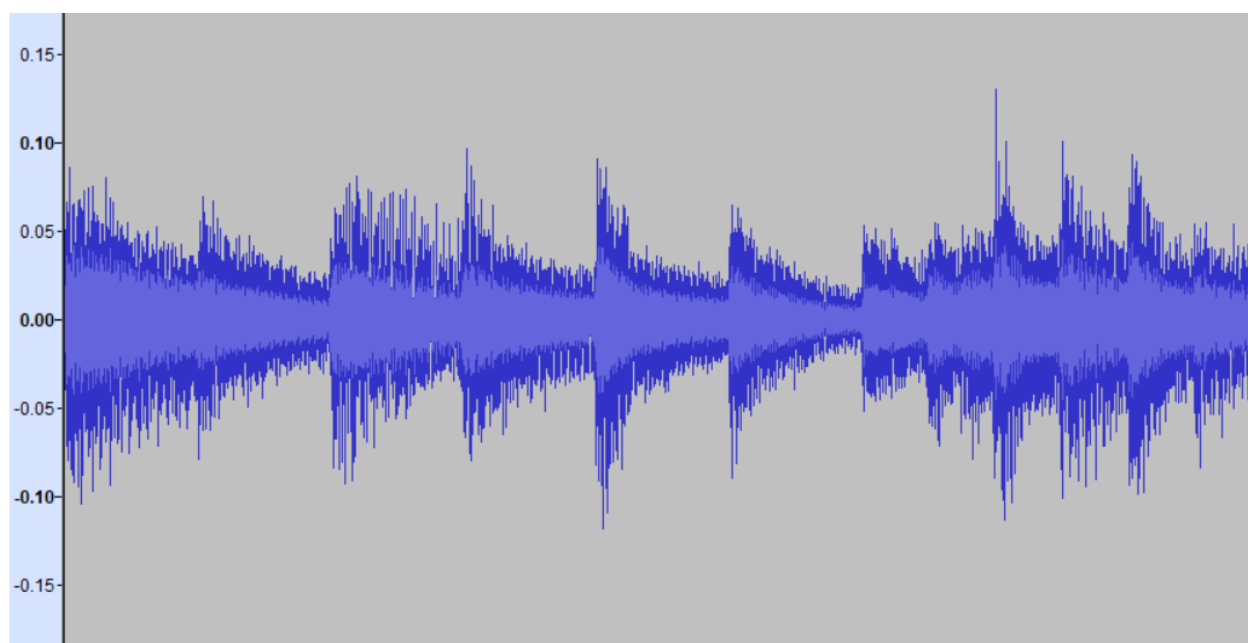
Rysunek 2.1.: Fragment chorału J.S. Bacha

Na pięciolinia umieszczone są takie elementy jak nuty, pauzy oznaczenia dynamiki, tempa i inne znaki dotyczące oznaczenia w jaki sposób należy wykonać dany utwór. Najpierw muzykę zapisywano na papierze, jednak po pojawieniu się komputerów zapis nutowy jak i samo nagranie odtworzenia mogło zostać zapisane na dysku. Nuty typowo są zapisywane jako obrazy lub pliki PDF, jednak dane zapisane w taki sposób nie nadają się do edycji lub przetwarzania przez algorytmy komputerowe. Jednym z pierwszych ustandaryzowanych formatów był *MusicXML*, który przy pomocy zwykłych tagów XML opisywał notację muzyczną.

Przykładowy zapis muzyki zapisano na obrazku 2.1. Przedstawiony fragment składa się z czterech głosów grających na raz oraz z czterech taktów.

2.2.1. WAV i MP3

WAV (ang. *Waveform audio format*) jest binarnym zapisem plików, który jest znany z tego, że jest w stanie zapisywać dźwięk nie używając żadnego algorytmu kompresji. W związku z tym rozmiary tych plików są bardzo duże, co sprawia, że ich przechowywanie wymaga znacznego miejsca na dysku co może być problematyczne, kiedy zbiór danych jest bardzo duży. Plik WAV jest najrzetelniejszą cyfrową reprezentacją dźwięku analogowego. Tak duży rozmiar danych zawdzięcza się temu, że dźwięk jest zapisywany z częstotliwością 44.1 kHz, czyli 44100 sampli na sekundę. Format został stworzony w roku 1991 roku, jest jednym z najbardziej rozpowszechnionych formatów i jest obsługiwany przez praktycznie każde



Rysunek 2.2.: Plik .wav otworzony w programie Audacity.

oprogramowanie edycji dźwięku. W celu wizualizacji dokonano syntezy?? utworu przedstawionego w postaci nutowej na obrazku 2.1 w związku z czym powstała fala dźwiękowa przedstawiona na grafice 2.2.

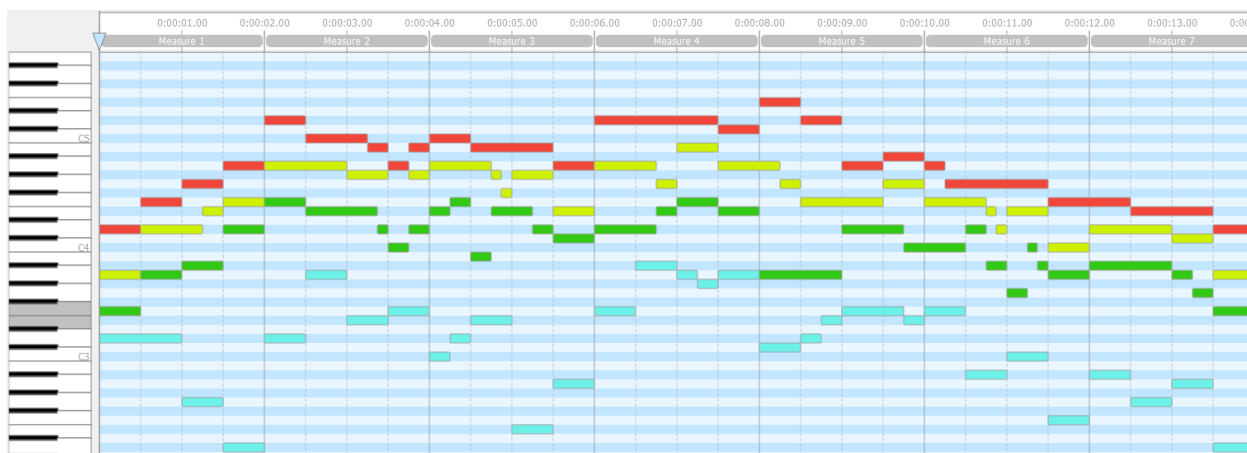
Często nie jest potrzebne przechowywanie bezstratne dźwięku, ponieważ i tak większość informacji, można zachować przy użyciu mniejszej ilości danych. Jednym z najpopularniejszych standardów kodowania stratnego jest MP3. Kompresja zmniejsza dokładność kodowania jak i również “ucina” częstotliwości, które nie są słyszalne dla człowieka. Stratne kodowanie próbuje zachować równowagę pomiędzy jakością dźwięku a rozmiarem pliku. W kontekście uczenia maszynowego, zmniejszona ilość sampli przy zachowaniu większości informacji jest zjawiskiem pożądanym, ponieważ udaje nam się przynajmniej częściowo rozwiązać problem związany z “przekleństwem wymiarowości”.

2.2.2. MIDI (ang. *Musical Instrument Digital Interface*)

MIDI jest standardem, który opisuje protokół komunikacji, interfejs cyfrowy oraz złącze pozwalające połączyć ze sobą elektroniczne instrumenty, komputery oraz inne muzyczne periferia. Pojedynczy kabel MIDI jest w stanie przekazać informacje na temat szesnastu na raz nadających kanałów z czego każdy może pochodzić z innego instrumentu. Każda interakcja z instrumentem, czyli przykładowo naciśnięty klawisz, szarpnięta struna, jest zapisana jako

“event”, który zachowuje wartości takie jak konkretny znacznik czasu, wysokość dźwięku czy jego głośność. Dane pochodzące z urządzeń są zapisywane w specjalnym pliku, głównie o rozszerzeniu *.mid* lub *.midi*. Plik pozwala na przechowanie, rozpowszechnianie jak i również edycje dźwięków. W związku z tym, że w pliku nie jest przechowywany zapis konkretnej fali dźwiękowej zapisanej przez mikrofon, pozwala to na późniejszą zmianę np. instrumentu, który będzie odtwarzał zapisane dźwięki.

Typowym przedstawieniem zapisanych danych jest tak zwany *piano roll*, który można porównać do dwuwymiarowego układu współrzędnych, gdzie współrzędna horyzontalna opisuje czas, a horyzontalna konkretne dźwięki pokazane jako klawisze pianina. Przykładowe porównanie zapisu MIDI z zapisem nutowym przedstawiono na grafice 2.3 oraz 2.1.



Rysunek 2.3.: Muzyka zapisana w pliku MIDI

Warto zaobserwować, że zapis pliku w formacie MIDI umożliwia zdecydowanie większą ekspresję, ponieważ jest to zapis odtworzenia przez artystę pewnego utworu. Ten sposób nie ogranicza muzyki sztywno do konkretnego rytmu zdefiniowanego przez autora.

2.2.3. Notacja ABC

Notacja ABC jest systemem zapisywania nut w postaci czystego tekstu znakami ASCII. Notacja pojawiła się w latach 70 XX wieku w celu zapisu i nauki tradycyjnych irlandzkich melodii, a następnie w kolejnej dekadzie została rozwinięta przez Chrisa Walshawę, który zapisywał w niej tradycyjne melodie zanim nauczył się standardowego zachodniego zapisu nutowego. Zapis ten posłużył do stworzonego przez niego programu *abc2mtex*, który na podstawie notacji ABC generował komendy pozwalające zapis partytur w postaci *MusicTex*. Obecnie używanym standardem jest wersja z roku 2011.

```

1 X:1
2 T:Chorał no.7
3 A:J.S. Bach
4 Q:1/4=120
5 V:1
6 L:1/16
7 M:4/4
8 K:C clef=G2
9 D4F4G4A4|d4c6B2A2B2|c4B8A4|d12^c4|
10 V:2
11 L:1/16
12 M:4/4
13 K:C clef=G2
14 A,4D6E2F4|A8^G4A2^G2|A6^G^F^G4E4|A6G2B4A4|
15 V:3
16 F,4A,4^A,4D4|F4E7DC2D2|E2F2B,2E4D2^C4|D6E2F4E4|
17 V:4
18 L:1/16
19 M:4/4
20 K:C clef=F4
21 D,8G,,4D,,4|D,4A,4E,4F,4|C,2D,2E,4E,,4A,,4|F,4^A,4A,2^G,2A,4|

```

Rysunek 2.4.: Zapis wielu głosów w notacji ABC.

Jak widać notacja ABC w dość zwięzły sposób zapisuje partyturę. Poza konkretnymi nutami oraz rytmem w tym formacie możemy zapisać również dodatkowe informacje na temat utworu. Każda linijka zaczynająca się od znaku A-Z a następnie dwukropkiem jest tak zwanym “polem informacyjnym”. W tych polach zapisywać można takie informacje jak tytuł utworu lub metrum oraz wiele innych informacji między innymi dotyczące z jakiego zbioru muzycznego pochodzi dany utwór. Wiele z tych informacji powinna być usunięta w procesie preprocessingu danych, tak aby model dostał tylko te informacje, które rzeczywiście pomogą w nauce struktury muzyki. Do takich pól należy przede wszystkim domyślna długość nuty (L:), metrum (M:) oraz tonacja (K:) w jakiej utwór został napisany. Każda z tych informacji jest możliwa do “odgadnięcia” przez model, jednak podając modelowi te informacje wyraźnie, mamy większą kontrolę nad procesem treningu. Dodatkową zaletą podawania tych informacji podczas treningu jest możliwość podania ich jako początkowa sekwencja na podstawie której model dalej będzie próbował kończyć melodię, przez co mamy kontrolę np. nad tym w jakiej tonacji i w jakim metrum zostanie wygenerowany nasz utwór.

Notacja ABC wspiera również melodie polifoniczne przy pomocy tagów V. Ich ilość nie jest ograniczona w związku z czym można w tym zapisie jest możliwe ujęcie nawet muzyki orkiestrowej. Dość skrajnym przykładem jest zapis drugiej części VII symfonii Ludwiga van Beethovena, która składa się aż z 19 instrumentów [2].

Tokenizacja plików MIDI

Token jest odrębnym elementem, częścią sekwencji tokenów. W języku naturalnym tokenem może być znak, zaimek lub słowo. Zdanie może być następnie tokenizowane na sekwencję tokenów reprezentujących słowa i znaki interpunkcyjne. W przypadku muzyki tokeny mogą reprezentować wartości atrybutów nut (wysokość, wartość, czas trwania) lub zdarzenia czasowe. Token może przyjmować jedną z trzech form:

- nazwa tokenu - słowna reprezentacja *eventu* MIDI np. *Pitch_50*
- id - unikalna wartość liczbowa przypisana konkretnemu zdarzeniu
- bajt - unikalny bajt, który został przydzielony podczas treningu tokenizera

Słownictwo jest zapisywane w postaci *look up table* łączącej nazwę tokenu z odpowiadającym jej id lub bajtem. Trening tokenizera polega na obliczeniu kodowania gramatykowego np. *byte pair encoding* do postaci tabelarycznej w celu wykorzystania ich w dalszym modelowaniu.

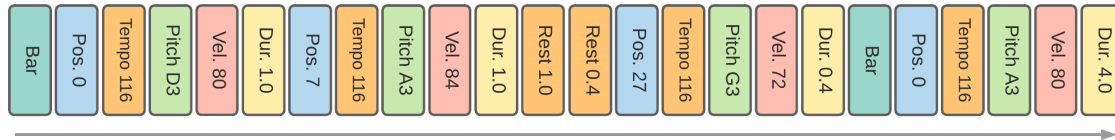
Poza tokenami, które tokenizer tworzy na podstawie pliku MIDI dodaje on również dodatkowe takie jak:

- PAD (*padding*) - token używany w przypadku kiedy w *batchu* długość sekwencji jest różna; w takim przypadku tym tokenem wydłuża się sekwencje aby wszystkie miały długość najdłuższej
- BOS (*beginning of section*) - token oznaczający początek sekwencji
- EOS (*end of section*) - token oznaczający koniec sekwencji

Istnieje wiele algorytmów tokenizacji MIDI, jednak aby przedstawić mechanizm działania, przyjrzymy się popularnemu algorytmowi REMI. REvamped MIDI (REMI) reprezentuje zdarzenia jako sekwencje tokenów wysokości tonu, dynamiki (głośności), długości trwania oraz czasu poprzez tokeny taktu i pozycji. Token taktu oznacza początek nowego taktu, natomiast token pozycji określa miejsce zdarzenia w bieżącym takcie. Porównanie pomiędzy zapisem nutowym a jego reprezentacją w formie tokenów można zaobserwować na ilustracjach 2.5 i 2.6. Tokeny **Bar** oznaczają początek taktu, **Pos** pozycję nuty w danym takcie, **Pitch** wysokość nuty, **Vel** jej głośność, a **Dur** czas jej trwania. Wiele algorytmów tokenizacji jest zaimplementowanych w bibliotece MidiTok dla języka Python [3].



Rysunek 2.5.: Prosta sekwencja muzyczna



Rysunek 2.6.: Reprezentacja muzyki w postaci tokenów

2.2.4. Porównanie zapisu muzycznego

W tabeli 2.1 zostało przedstawione porównanie zapisu całego utworu, którego fragment przedstawiono na grafice 2.1.

	rozmiar pliku na dysku	długość sekwencji
wav	3971116 B	992768
mp3	360951 B	992768
midi	1638 B	740
abc	953 B	562

Tabela 2.1.: Porównanie różnych zapisów muzyki

Jak widać zapis w postaci fali dźwiękowej zajmuje najwięcej miejsca na dysku, jak i również długość sekwencji, która zapisuje cały utwór jest najdłuższa. W celu uzyskania sekwencji z pliku MIDI został użyty tokenizer REMI o którym jest mowa w rozdziale 2.2.3. Przy użyciu innego tokenizera niż REMI, liczba powstałych tokenów może lekko różnić się od tej w tabeli. Różnica w rozmiarze na dysku oraz długością sekwencji pomiędzy zapisem MIDI oraz ABC nie jest widoczna, jednak oba te formaty są zdecydowanie lepszymi rozwiązaniami w przypadku generacji notacji muzycznej (dziwne zdanie do poprawy raczej). Warto zwrócić uwagę, że gdyby w oryginalnym utworze pojawiły się powtórzenia danych sekwencji melodii notacja ABC byłaby jeszcze bardziej oszczędna w obu przypadkach od plików MIDI. Powodem tego jest fakt, że notacja ABC zapisuje bezpośrednio w postaci tekstowej zapis nutowy, czyli znak powtórzenia zostanie zapisany jako :| oraz |: , a w pliku MIDI cała powtarzana sekwencja musi zostać ponownie odtworzona przez artystę.

Do głównej zalety plików MIDI należy ich wszechstronność. Nie są one ograniczone do kodowania sztywno określonego rytmu tak jak w notacji muzycznej. Dodatkowo jest bardzo

łatwo przy użyciu narzędzi komputerowych zamienić dowolny zapis nutowy (MusicXML, ABC) na plik MIDI. Niestety w drugą stronę nie jest już tak łatwo, właśnie w związku z nie konkretnym rytmem, który może zostać zakodowany w MIDI. Istnieją narzędzia takie jak np. MuseScore, które jest w stanie wygenerować partyturę z pliku MIDI, jednak nie jest ono w każdym przypadku jednakowo skuteczne.

Plik w postaci fali dźwiękowej nie jest idealnym sposobem zapisu muzyki, jednak jest to często jedyne rozwiązanie aby zebrać dużą ilość danych. Szczególnie jeśli mamy do czynienia z muzyką nowoczesną zazwyczaj twórcy nie udostępniają plików MIDI lub zapisu nutowego na podstawie którego powstała dana piosenka. Z tego powodu zostaje albo praca z plikami mp3 lub czasochłonna translacja muzyki na jej zapis w innym formacie. Dodatkowo zaletą plików z dźwiękiem w zapisie cyfrowym jest dostęp do potencjalnie istniejącego zapisu odśpiewanego tekstu. Co prawda przykładowo zapis w notacji ABC pozwala na przekazanie tekstu w tagu W:, jednak jest to tylko jego zapis a nie faktyczne odśpiewanie. Z tego powodu pliki wav nadają się np. do próby replikacji czyjegoś głosu.

Główną zaletą plików ABC jest ich tekstowa reprezentacja notacji muzycznej. Z tego powodu nie wymagają one wiedzy na temat struktury, budowy oraz standardów plików np. MIDI aby zacząć z nimi pracować. Ponieważ notacja ABC jest tekstowa, można ją łatwo integrować z narzędziami do przetwarzania tekstu. Zwięzłość zapisu jest również przyczyną dla której model uczenia maszynowego może być mniejszy w związku z krótszą ilością tokenów w sekwencji, na podstawie której model się uczy. Zapis w postaci tekstu pozwala na bardzo proste dodanie pewnych dodatkowych danych, które można policzyć w etapie *preprocessingu* danych. Podejście takie zostało zaproponowane w modelu *Tunesformer*[4], który wziął przykład z artykułu CTRL [5], w którym autorzy dodawali do korpusu treningowego takie tagi jak *Books*, *Horror*, *Relationships*, *Legal*. Tagi te były podawane na początku *promptu* podawanego dla modelu, przez co model "wiedział" w jakim stylu powinien udzielić odpowiedzi. W przypadku muzyki w zapisie ABC zostały wprowadzone nowe *control codes*:

- S: - liczba sekcji w całym utworze. Tag jest prost do policzenia ponieważ początek i koniec taktu jest jawnie oznaczany przy pomocy symboli [l, ll, ll], |: , :: i :|.
- B: - liczba taktów w danej sekcji. Liczy tylko wystąpienia znaku l
- E: - podobieństwo. Kontroluje poziom podobieństwa pomiędzy dwoma następującymi po sobie sekcjami. W związku z tym, że zapis melodii jest w postaci tekstu, można policzyć podobieństwo używając odległości Levenshteina używając wzoru

$$E(c, p) = 1 - \frac{lev(c, p)}{\max(|c|, |p|)} \quad (2.2)$$

gdzie c oraz p to następujące po sobie sekwencje, a $|c|$, $|p|$ to ich długości.

Policzone kody wraz z ich odpowiednim oznaczeniem dodawane są do pliku tekstowego, który zapisuje notacje. Dodanie ich nie powoduje, że notacja staje się nieprawidłowa. Tak samo jak w przypadku tagów, które opisują elementy zapisu jak tonacja czy metrum, służą one nam oraz modelowi jako wskazówki, którymi powinien się sugerować podczas generacji muzyki. Wartości podawane w tagach są elementami, które policzyliśmy na podstawie danych, co znaczy że model mógłby je w czasie treningu wywnioskować jednak podawanie ich w sposób *explicite* przekazuje użytkownikowi kontrolę nad tymi parametrami.

2.3. Zbiory danych

2.3.1. Johann Sebastian Bach Chorales

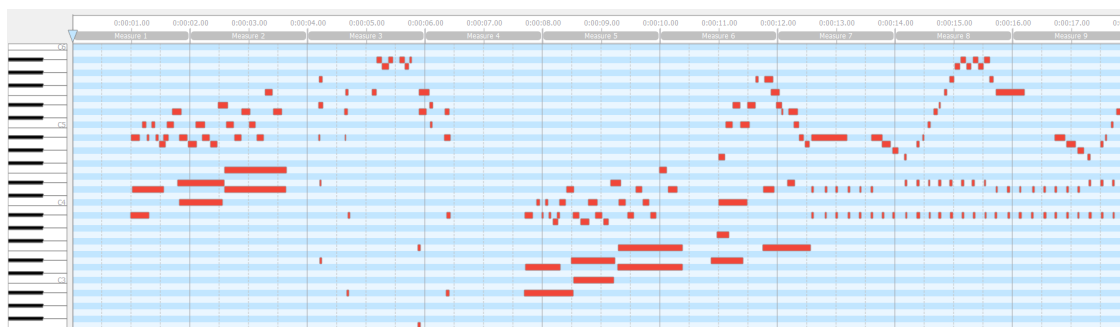
JSB Chorales [6] to zbiór krótkich, czterogłosowych utworów. Chorały zostały pierwotnie skomponowane przez Johanna Sebastiana Bacha w XVIII wieku. Napisał je najpierw biorąc wcześniej istniejące melodie ze współczesnych hymnów luterkańskich, a następnie harmonizując je w celu stworzenia części dla pozostałych trzech głosów. Wersja zbioru danych używana kanonicznie w kontekście uczenia się reprezentacji w kontekście uczenia się reprezentacji składa się z 382 takich chorałów. Utwory przed złożeniem w zbiór danych zostały ztransponowane do tonacji C-dur przez co stają się poostrze do analizy przez algorytmy. Dodatkowo ich rytm został znormalizowany, i najmniejszą wartością jest szesnastka. Dzięki temu uproszczeniu zbiór jest bardzo łatwo konwertowalny do dowolnego zapisu które zostały zaprezentowane na grafikach 2.1 2.3 2.4.

Dodatkową zaletą takiego rytmu jest możliwość pominięcia użycia tokenizera MIDI, ponieważ wiedząc że rytm jest stały i wyznaczony przez szesnastkę, można użyć naiwnego podejścia i zapisać każdy głos jako wektor, którego wartościami jest obecnie grająca nuta. Po złożeniu wektorów w macierz $4 \times$ długość utworu można przejść do analizy. Wadą tego podejścia jest możliwość wystąpienia macierzy rzadkich, które często stanowią problem dla algorytmów uczenia maszynowego (CITATION NEEDED!), jednak w tym zbiorze ten problem nie występuje.

2.3.2. The MAESTRO v3.0

Zbiór danych MAESTRO w wersji trzeciej [7] powstał w współpracy z *Minnesota International Piano-e-Competition*, czyli międzynarodowym konkursem pianistycznym. Litera *e* w nazwie odnosi się do fortepianów używanych podczas konkursu którymi są Yamaha Disklavier. Jest to rodzaj fortepianu, który poza tradycyjną funkcjonalnością jest obudowany elektronicznymi sensorami, które pozwalają między innymi zapis odegranej muzyki w formie MIDI. Z tego powodu organizacja Magenta działająca wewnątrz Google użyła zapisu melodii z wielu edycji konkursu aby stworzyć dataset zawierający około 200 godzin muzyki

fortepianowej. W tych 200 godzinach znajduje się 1276 występów w skład których wchodzi ponad 7 milionów nut. Dwa rodzaje zbioru jakie można pobrać to zbiór w wersji WAV oraz MIDI. Pierwszy, z powodów, które przedstawiono w tabeli 2.1 zajmuje aż 122 GB, natomiast wersja MIDI zajmuje tylko 81 MB. Jak widać zapis plików WAV jest zdecydowanie mniej oszczędny jeśli chodzi o przechowywanie danych na dysku.

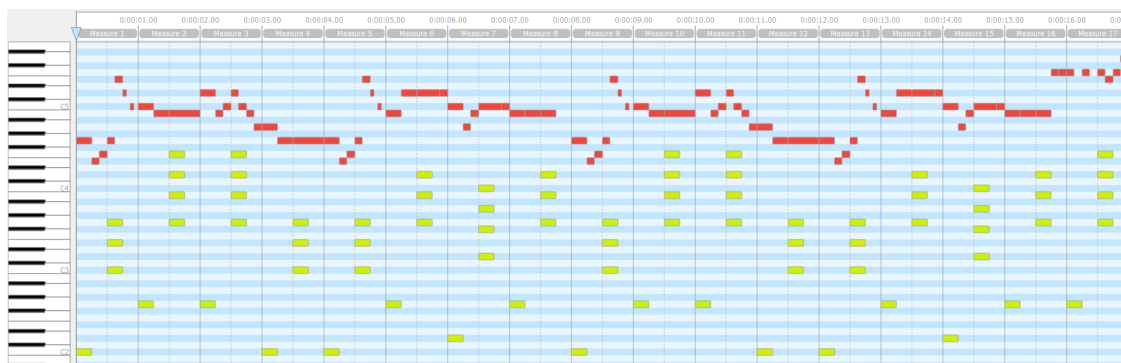


Rysunek 2.7.: Przykład utworu pochodzącego ze zbioru MAESTRO.

Porównując zapis prostego utworu widoczny na grafice 2.3 z tym na grafice 2.7 osoby nawet nie znające się na muzyce są w stanie wywnioskować że ten utwór bardziej złożony. Jak widać odegrane nuty często nie padają na sztywno określone rozpoczęcie taktu, przez co potencjalnie wygenerowana muzyka na podstawie tego zbioru będzie wydawała się bardziej ekspresyjna oraz “ludzka”.

2.3.3. IrishMAN

Zbiór danych IrishMAN (*ang. Irish Massive ABC Notation*) [8] zawiera 216,248 fragmentów melodii podzielonych na zbiór treningowy i walidacyjny. Melodie zebrane są ze stron takich jak thesession.org oraz abcnotation.org, które znane są z udostępniania tradycyjnej muzyki. Aby zapewnić jednolitość formatowania najpierw wszystkie utwory zostały przekonwertowane do notacji MusicXML, a następnie do notacji ABC. Istnieje bliźniacze zbiory IrishMAN-MIDI oraz IrishMAN-XML, które zawierają te same utwory jednak zapisane w innych formatach. Wszystkie fragmenty muzyki należą do domeny publicznej, w związku z czym bez obaw o łamanie praw autorskich można korzystać z tych melodii. Zbiór jest o tyle ciekawy że zawiera w sobie muzykę jedno jak i wielogłosową. Autorzy zbioru udostępnili również skrypty napisane w języku Python przy pomocy których, można samemu stworzyć własne zbiory danych na podstawie swoich plików ABC lub MusicXML. Jest to również jeden z niewielu zbiorów, który udostępnia te same fragmenty muzyki w różnych formatach, przez co nadaje się on wyjątkowo dobrze do porównywania różnych algorytmów generacyjnych.

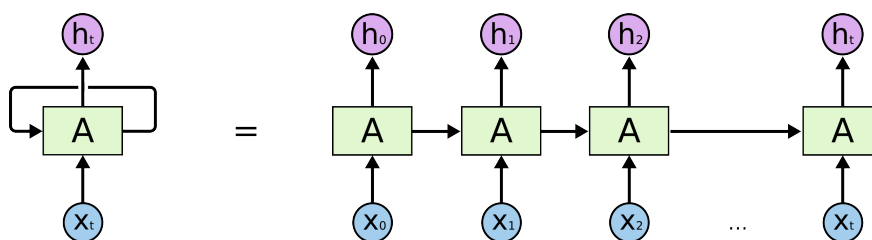


Rysunek 2.8.: Przykład wielogłosowego fragmentu ze zbioru IrishMAN.

2.4. Architektury transformera

2.4.1. Sieci RNN

Pierwszymi próbami użycia sieci neuronowych w przetwarzaniu sekwencji było użycie rekurencyjnych sieci neuronowych. Architektura ich była prosta i dodawała ona pętlę w taki sposób aby wartości mogły być przekazywane pomiędzy krokami. Rysunek 2.9 przedstawia fragment sieci oraz rozwinięcie pętli.



Rysunek 2.9.: Schemat budowy fragmentu sieci RNN.

Prezentowane podejście jest dość proste, jednak niesie ze sobą pewne problemy.

Pierwszym z nich jest brak selekcji, które informacje są rzeczywiście istotne w danym kontekście, a które nie. W zadaniu przewidywania kolejnego słowa np. dla zdania “Kupiłem zegarek na *rękę*” nie jest wymagany żaden dodatkowy kontekst, a odległość pomiędzy ważnymi informacjami jest mała. Jednak dla zdania “Rok temu ukończyłem kurs lotnika ... Teraz staram się o licencje *pilota*” wymagany jest szersze spojrzenie na całą strukturę wypowiedzi i “przypomnienie” informacji z początku zdania. W przypadku muzyki sytuacja ma się dokładnie tak samo. Przykładowo, aby rozwiązać akord dysonansowy, wymagana jest

jedynie informacja na jego temat, tak aby dobrać odpowiedni akord konsonansowy lub akord dominantowy na tonikę, jednak jeśli w muzyce przewija się pewien temat, to wymagane jest zachowanie informacji w sieci, w jakich momentach się on pojawia, aby móc dodać go w odpowiednie miejsce. W teorii zwyczajna sieć rekurencyjna jest zdolna do zapamiętywania długich zależności, jednak w praktyce nie jest to bardzo rzadko osiągalne.

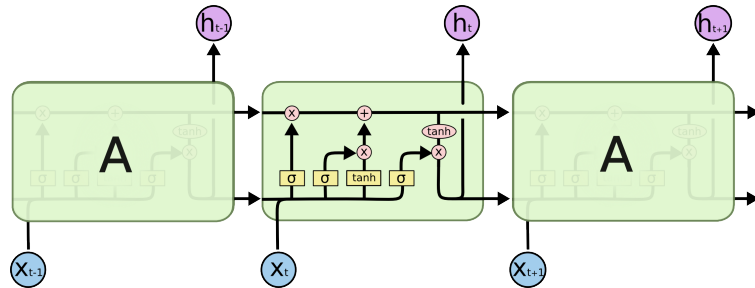
Dodatkowym problemem podczas treningu modeli jest problem zanikającego gradientu.

$$\begin{aligned}
 H_{i+1} &= A(H_i, x_i) \\
 H_3 &= A(A(A(H_0, x_0), x_1), x_2) \\
 &\text{gdzie:} \\
 A(H, x) &:= Wx + ZH \\
 H_N &= W^N x_0 + W^{N-1} x_1 + \dots
 \end{aligned} \tag{2.3}$$

Jak widać z powyższego równania, kiedy N będzie wystarczająco duże (a na pewno będzie ponieważ jest to długość okna kontekstowego) macierz wag W będzie podnoszona do bardzo wysokiej potęgi przez co jej wartości poniżej 1 będą zanikać do 0, a te będące powyżej 1 będą eksplodować w nieskończoność.

2.4.2. LSTM

Aby rozwiązać omawiane problemy sieci rekurencyjnych została zaproponowana architektura sieci nazwana *Long Short Term Memory*. Ich budowa pozwala na długofalowe zapamiętywanie i przekazywanie informacji. Zazwyczaj są łączone między sobą co pozytywnie wpływa na działanie sieci.



Rysunek 2.10.: Komórki LSTM połączone między sobą.

Główną częścią budowy komórki jest górna strzałka na obrazku 2.10, która posiadając jedynie proste operacje liniowe przekazuje z stany ukryte z jednej części do drugiej. Dolna część komórki zawiera tak zwane bramki, które kontrolują jaka które informacje są ważne i mają być zachowane a które można zignorować. Są one zbudowane z warstw sigmoidalnych oznaczonych na grafice literami σ . Wartościami zwrótnymi jest liczba pomiędzy 0 a 1 gdzie

mała wartość znaczy niską wagę, a wysoka dużą. LSTM zwraca dwa wektory. Pierwszy to stan ukryty (*hidden state*), który przechowuje informacje o bieżącym stanie sieci i jest aktualizowany przy każdym kroku czasowym. To ten wektor jest bezpośrednio wykorzystywany do przewidywania kolejnych słów lub innych zadań wyjściowych. Drugi nazywany stanem komórki (*cell state*) jest dodatkowy i przechowuje informacje na dłuższą metę, przez co jest odpowiedzialny za przechowywanie bardziej trwałego kontekstu. Stan komórki jest modyfikowany przez tzw. bramki (*gate*), które decydują, które informacje powinny być dodane, usunięte lub zaktualizowane.

Architektura takiego modelu sprawia, że jest on o wiele mniej podatny na problem zanikającego gradientu, jednak jest to sytuacja, która niestety może nadal występować. Dodatkową wadą modelu LSTM jak i klasycznego RNN jest ich wysoka złożoność treningu, która występuje z powodu konieczności użycia specjalnego wariantu algorytmu wstecznej propagacji błędów nazywanego “wsteczną propagacją błędów w czasie” *backpropagation through time*, która rozwija sieci rekurencyjne, aby propagować gradienty.

Istnieje kilka wariantów LSTM-ów np. dwukierunkowe LSTM, które pozwalają aby informacja przepływa w dwóch kierunkach na raz a nie tylko z lewej strony do prawej. Użycie takich modeli poprawia skuteczność i wyniki w stosunku do klasycznego modelu, jednak problemy z ich treningiem dalej pozostają lub nawet stają się jeszcze większe.

Uwaga 2.1. Można coś więcej opisać ale czy jest po co? Jak będzie brakowało stron a będzie czas to się dopisze _(\`)/ Plus trzeba dodać cytowania

2.4.3. Algorytm uwagi (ang. *attention*)

Mechanizm uwagi (ang. *attention*) został zaproponowany w artykule “Attention is All You Need” [1]. Algorytm ten działa w ustalonym oknie kontekstowym (ang. *context window*), w którym przydziela każdemu słowu wagę reprezentującą jego wagę w danym kontekście. Metoda ta imituje sposób, w jaki człowiek skupia się na różnych fragmentach tekstu w danym momencie. Mechanizm uwagi został opracowany, aby rozwiązać problem występujący w rekurencyjnych sieciach neuronowych (*RNN*), gdzie “waga” słowa była tym większa, im bliżej końca sekwencji się znajdowało. W rezultacie modele często ignorowały informacje znajdujące się na początku okna kontekstowego.

W kontekście muzyki problem zanikania kontekstu jest jeszcze bardziej widoczny niż w tradycyjnych problemach przetwarzania języka naturalnego (NLP), ponieważ na początku zapisu nutowego znajdują się kluczowe informacje, takie jak metrum czy tonacja utworu, zapisane za pomocą krzyżyków lub bemoli. Z tego powodu model, analizując dalsze części sekwencji, może ignorować te istotne elementy.

Zależności pomiędzy słowami w mechanizmie uwagi można przedstawić jako graf pełny skierowany, gdzie wierzchołkami są poszczególne słowa lub inne tokeny z okna czasowego, a

krawędziami są wartości określające, jak istotne są dla danego słowa wszystkie inne. Jest to zdecydowanie lepsze podejście niż w rekurencyjnych sieciach neuronowych lub LSTM-ach, gdzie zależności dla danego słowa są przekazywane jedynie w wektorze, który kompresuje kontekst ze wszystkich poprzednich słów. Dzięki temu mechanizm uwagi pozwala na bardziej dynamiczną kontrolę i monitorowanie, jak zależności pomiędzy tokenami formują się podczas treningu.

Uwaga 2.2. Tutaj trzeba napisać więcej na temat attention + jakiś obrazek jak to działa i opisać czym są K , V i Q

2.4.4. Transformer

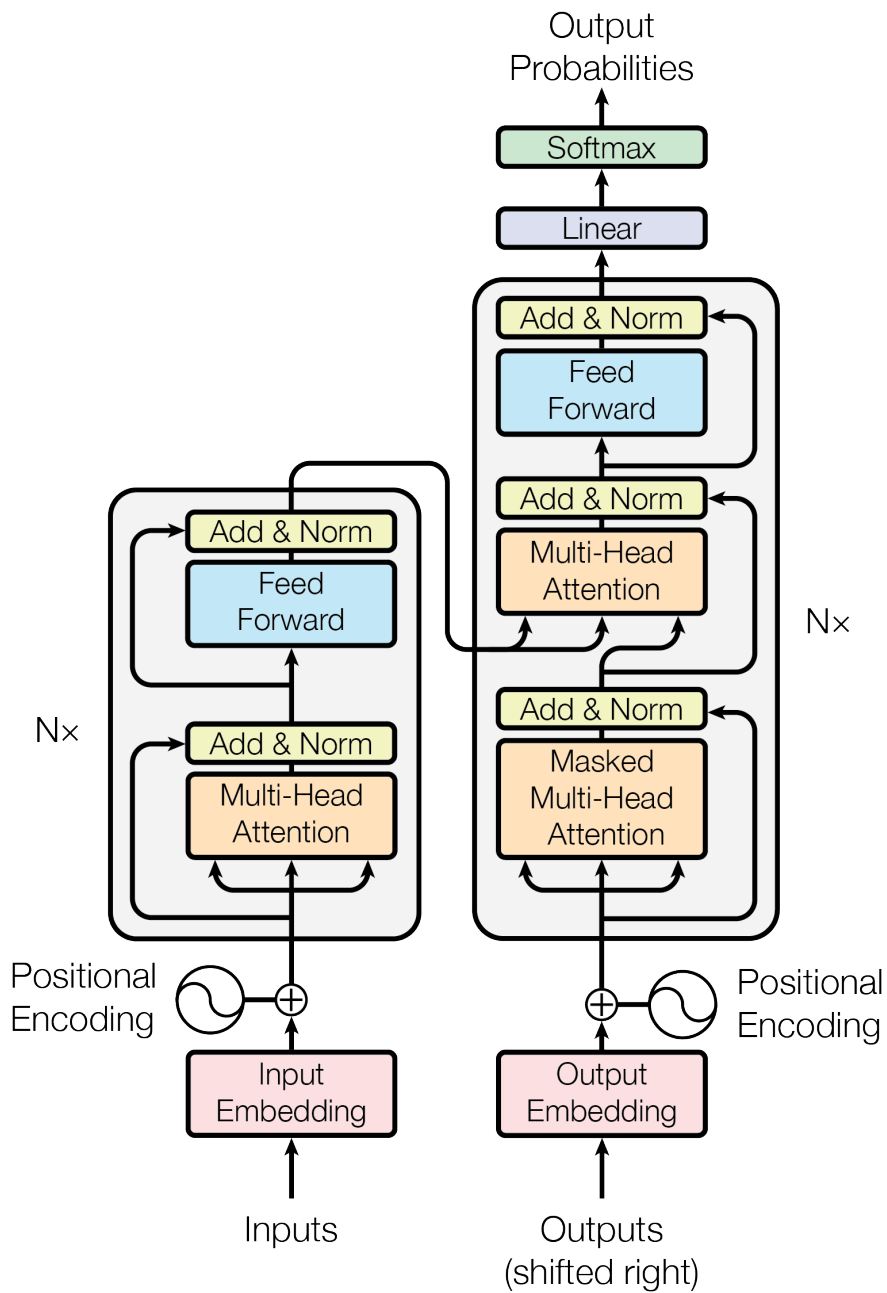
Wraz z wprowadzeniem algorytmu uwagi pisanego w 2.4.3 w tym samym artykule został zaproponowana architektura modelu, która wykorzystywała *multi-headed* algorytm uwagi. Cały model składa się z enkodera oraz dekodera.

Enkoder

Zadaniem enkodera jest zapoznanie się z całym korpusem tekstu i przekazanie zebranych informacji do dekodera. Mechanizm samo-uwagi *self-attention* jest centralnym elementem każdej warstwy enkodera. Pozwala on modelowi analizować relacje między różnymi tokenami w sekwencji niezależnie od ich odległości od siebie. Drugi komponent każdej warstwy enkodera to sieć neuronowa typu *feed forward*, która składa się z dwóch w pełni połączonych warstw z funkcjami aktywacji *ReLU*. Każdy z dwóch głównych komponentów jest otoczony mechanizmem dodania rezydualnego (poprawiającego wydajność treningu) i warstwą normalizacji. Dane przed wejściem są kodowane przy pomocy tabeli kodowań do wektorów oraz, aby zapewnić modelowi informacje na temat kolejności słów we fragmencie, jest dodawany wektor kodujący ich kolejności.

Dekoder

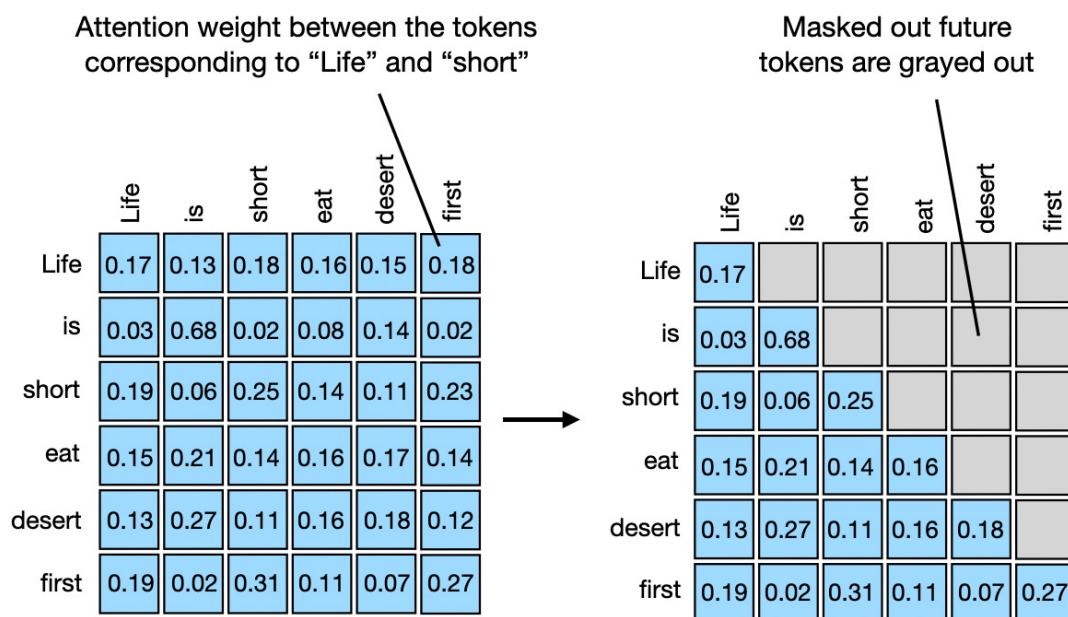
Zadaniem Dekoder w architekturze transformerowej jest generowanie wyjścia na podstawie sekwencji wejściowej i dotychczas wygenerowanych tokenów wyjściowych. Podobnie jak w enkoderze, mechanizm samo-uwagi analizuje zależności między tokenami w sekwencji wyjściowej, jednak w dekodерze stosuje się maskowanie (*ang. masking*), aby zablokować dostęp do przyszłych tokenów podczas generowania bieżącego tokena. Maskowanie uniemożliwia modelowi zobaczenie tokenów, które są generowane w późniejszych krokach. Maskę można zilustrować jako macierz trójkątną, co zostało pokazane na grafice 2.12. Drugi mechanizm uwagi umożliwia dekodерowi skupienie się na odpowiednich częściach sekwencji wejściowej



Rysunek 2.11.: Schemat transformera.

(wyjścia enkodera) podczas generowania każdego tokena wyjściowego. Proces ten jest podobny do mechanizmu samo-uwagi, ale tutaj zapytania (Q) pochodzą z poprzedniej warstwy

uwagi dekodera, a klucze (K) i wartości (V) pochodzą z wyjść enkodera. Następnie tak jak w poprzednim przypadku występuje warstwa *feed forward* i warstwy normalizujące. W ostatnich warstwach znajduje się warstwa liniowa, której zadaniem jest wyprodukowanie wyjścia o rozmiarze ilości wszystkich słów lub tokenów, oraz warstwa *softmax*, która przypisuje każdemu tokenowi prawdopodobieństwo bycia następnym słowem w sekwencji.



Rysunek 2.12.: Algorytm uwagi z maską.

Istnieją modele transformerowe nie korzystające z całej architektury. Przykładem *encoder only transformer* jest model BERT [9], który jest typowo wykorzystywany do zadań np. klasyfikacji. W przypadku takiego zadania, nie potrzeba omawianej wcześniej maski, ponieważ model powinien analizować kontekst całego fragmentu w "obie strony". Wystarczy dołożyć jako ostatnią warstwę warstwę liniową o odpowiednim rozmiarze, oraz zastosować odpowiednią funkcję aktywacji, aby przerobić enkoder na model klasyfikacyjny. Analogicznie, jeśli celem zadania jest wyłącznie generacja tekstu lub muzyki, enkoder nie jest wymagany. W takim problemie należy usunąć drugą warstwę uwagi (bez maski), ponieważ wektory kluczy i wartości pochodzące z enkodera będą nieznane. Generacja będzie odbywała się jedynie na podstawie podanej sekwencji lub tokenu BOS (*begin of sequence*), co da modelowi największą wolność twórczą. Rozwiązania *decoder only* są obecnie bardzo popularne i są używane przez modele takie jak GPT od OpenAI. Ciekawą obserwacją jest fakt, że dekodery pozbawione drugiej warstwy uwagi, stają się w zasadzie enkoderem z dodaną maską w pierwszej warstwie.

Jest to często wykorzystywana “sztuczka”, którą można wykorzystać w wielu bibliotekach implementujących gotowe części składowe modeli.

Trening transformera jest możliwy do bardzo dobrego zrównoleglenia. Wynika to z faktu, że mechanizm uwagi jest w stanie niezależnie liczyć wartości pomiędzy słowami. Przykładowo w zadaniu “*Ala ma kota.*” wartości pomiędzy *ma* i *kota* można policzyć nie znając wartości pomiędzy *Ala* i *ma*. Niestety podczas generacji kolejnych tokenów na podstawie istniejącej sekwencji należy policzyć ponownie wartości uwagi, przez co dla sekwencji o długości n należy wykonać n^2 operacji. Jest to niestety problematyczne jeśli chcemy rozważać bardzo długie sekwencje. Omawiany wcześniej model RNN, podczas predykcji, nie ma problemów ze złożonością, ponieważ cała informacja na temat wcześniejszej sekwencji jest skompresowana w stanie ukrytym. W takim przypadku okno kontekstu jest teoretycznie nieskończone.

GPT-2

[10]

2.5. Architektura *state space*

2.5.1. Mamba

Modele przestrzeni stanów (*State Space Models* - SSM) to matematyczne modele dynamicznych systemów, które są szeroko stosowane w inżynierii, ekonomii, biologii i wielu innych dziedzinach. Modele te są szczególnie przydatne w analizie i prognozowaniu czasowych szeregów danych.

Typowo, dla każdej jednostki czasu t model:

1. mapuje sekwencje wejściową $x(t)$
2. do reprezentacji ukrytej $h(t)$
3. i przewiduje sekwencję wynikową $y(t)$

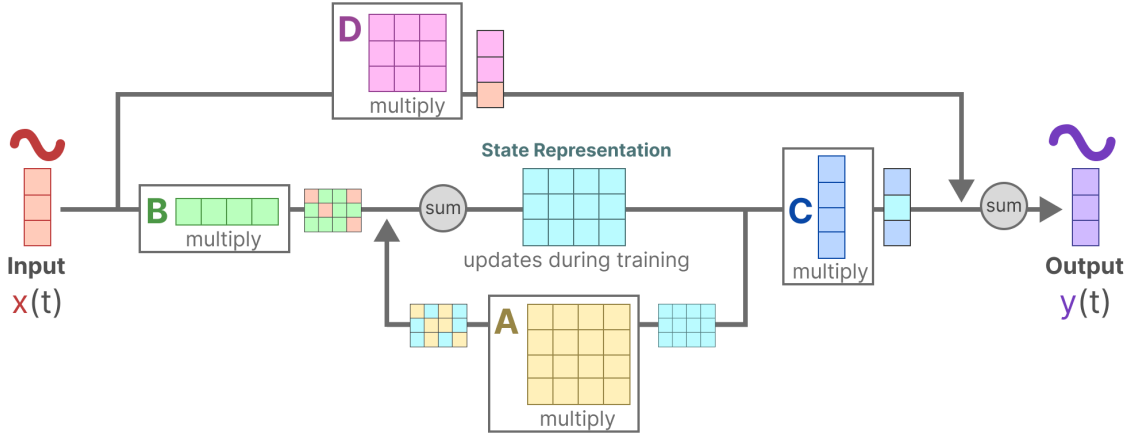
Modele SSM zakładają, że każdy system dynamiczny może zostać opisany w momencie czasu t używając dwóch równań.

$$\text{równanie stanu: } h'(t) = Ah(t) + Bx(t) \quad (2.4)$$

$$\text{równanie wynikowe: } y'(t) = Ch(t) + Dx(t) \quad (2.5)$$

Celem jest poznanie reprezentacji stanu $h(t)$ w taki sposób aby przejść z sekwencji wejściowej do sekwencji wynikowej.

Równanie stanu 2.4 opisuje w jaki sposób zmienia się stan w zależności od tego, jak dane wejściowe mają na niego wpływ. Wykonywane jest to przez macierze A oraz B . Celem równania 2.5 jest opisanie jak obecny stan ukryty wpływa na sekwencje wynikowym oraz jak sekwencja wejściowa wpływa na wynik. Te informacje są zapisywane odpowiednio w macierzy C oraz D . Schematyczne przedstawienie działania modelu zostało przedstawione na grafice 2.13.



Rysunek 2.13.: Schemat modelu przestrzeni stanów.

W związku z tym, że macierz D przypomina połączenia rezydualne lub połączenia typu *skip-connection* w głębokich sieciach neuronowych, często pomija się je w zapisie modelu co upraszcza zapis [11].

Dyskretyzacja modelu

Modele tego typu działają domyślnie w domenie danych ciągłych. Tekst lub muzyka są przykładami danych dyskretnych, przez co model musi zostać pod niego dostosowany. Aby dokonać dyskretyzacji modelu wykorzystuje się metodę eksploratora (lub interpolatora) rzędu zerowego (*ang. zero-order hold*). Metoda ta polega na podtrzymywaniu każdego otrzymanego sygnału, dopóki nie zostanie zapisany kolejny. Czas trwania poprzedniego sygnału jest dodatkowym parametrem nazywanym krokiem. W związku z tym, że model na wyjściu zwraca sygnał ciągły, aby zamienić go na dyskretny, należy *samplewać* go, z odpowiednią długością kroku. Po dyskretyzacji modelu można nazwać go modelem *sequence to sequence*. Aby lepiej zobrazować dokonane zmiany należy spojrzeć na równania 2.6 oraz 2.7.

$$\text{równanie stanu: } h_k = Ah_{k-1} + Bx_k \quad (2.6)$$

$$\text{równanie wynikowe: } y_k = Ch_k \quad (2.7)$$

Model operuje teraz na sekwencjach a nie na sygnałach. Notacja została zmieniona z t na k aby lepiej zobrazować, że obecnie model pracuje na dyskretnych znacznikach czasu. Dodatkowo w zapisie równania usunięto macierz D .

Reprezentacja konwolucyjna

Patrząc na zapis obu równań od razu można zaobserwować ich podobieństwo do rekurencyjnych sieci neuronowych. Każdy stan ukryty h zależy w głównej mierze od stanu z poprzedniego kroku czasowego. *SSM-y* mogą również być reprezentowane w postaci konwolucji. W związku z tym, że tekst lub muzyka są sekwencjami jednowymiarowymi w przeciwieństwie do typowego zastosowania sieci konwolucyjnych czyli obrazów, które są dwu lub trzy wymiarowe należy zastosować kernel o wymiarze 1. Tego rodzaju interpretacja modelu jest o tyle wygodna, że pozwala zrównoleglić trening modelu a co za tym idzie zmniejszyć jego czas [12]. Kiedy model jest już wytrenowany i gotowy do wykonywania generacji, można bez problemu przejść na interpretację rekurencyjną.

Macierz HiPPO

Macierz A modelu jest jego jednym z najistotniejszych elementów. To od niej zależy, czy model zapamięta cały kontekst wypowiedzi czy tylko kilka ostatnich tokenów. Jest to szczególnie ważne w przypadku predykcji, ponieważ model bazuje wyłącznie na podstawie ostatniego stanu ukrytego. Aby zapewnić dobrą kompresję kontekstu stosuje się macierz HiPPO [13]. Zadaniem *High-order Polynomial Projection Operators* jest jak najlepsze skompresowanie wejścia do wektora cech. Konstrukcje macierzy pokazano na równaniu 2.8. W tym przypadku zakładamy, że macierz jest kwadratowa.

$$A_{nk} = - \begin{cases} (2n+1)^{1/2}(2k+1)^{1/2} & \text{dla } n > k \\ n+1 & \text{dla } n = k \\ 0 & \text{dla } n < k \end{cases} \quad (2.8)$$

Zbudowanie macierzy w taki sposób prowadzi do zdecydowanie lepszych wyników niż zainicjalizowanie jej losowo. Macierz HiPPO reprodukuje sekwencje w taki sposób, że dane które są bliżej końca sekwencji są reprodukowane lepiej, niż te na początku. Ideą stojącą za używaniem takiego rozwiązania jest chęć zachowania jak największej liczby informacji w stanie ukrytym. Użycie współczynników wielomianów Lagrange'a pozwala aby macierz pomagała zapamiętywać długie zależności pomiędzy tokenami [14].

Modele 4S

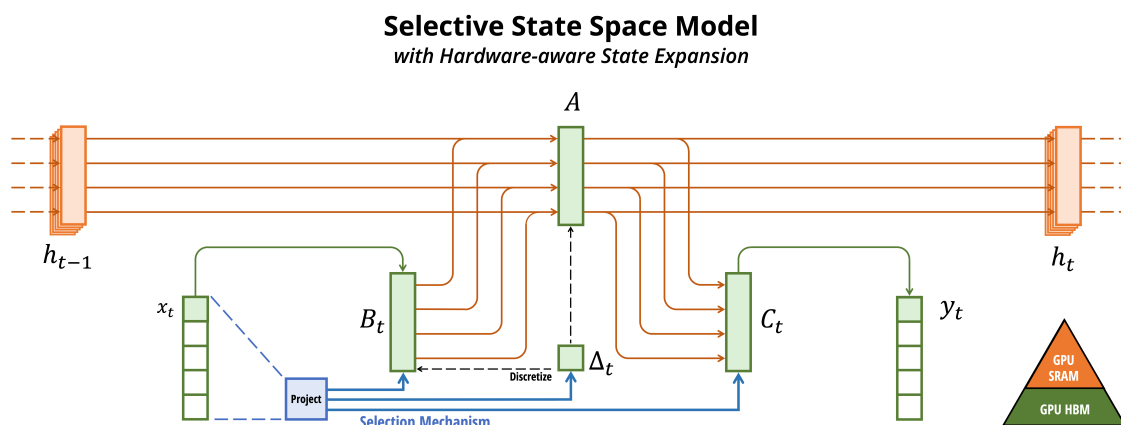
Po dodanie omawianych usprawnień (dyskretyzacja, macierz HiPPO) do tradycyjnego modelu SSM nowo otrzymany model dostał nazwę *Structured State Space for Sequences*

(4S). Tego rodzaju modele pozwalają na śledzenie długich zależności w sekwencji dzięki macierzy HiPPO, a możliwość ich reprezentacji w postaci rekurencyjnej oraz konwolucyjnej rozwiązuje problemy między innymi złożoności obliczeniowej treningu i infrencji modelu.

2.5.2. Mamba

Mamba [15] jest architekturą modelu, która bazuje na modelach 4S, jednak wprowadza dodatkowe usprawnienia takie jak:

1. algorytm selektywnego skanowania (*selective scan*) pozwala skutecznie skupiać się na istotnych danych wejściowych i ignorować te nieistotne. Zdolność ta umożliwia lepszą obsługę zależności dalekiego zasięgu i filtrowanie szumów, co stanowi znaczną poprawę w stosunku do standardowych transformatorów, które traktują wszystkie dane wejściowe jako tak samo ważne
2. algorytm, *hardware-aware*, który pozwala uniknąć konieczności wykonywania rozległych operacji wejścia/wyjścia pomiędzy różnymi poziomami hierarchii pamięci GPU, co prowadzi do znacznej poprawy szybkości. W szczególności, skaluje się ona liniowo wraz z długością sekwencji



Rysunek 2.14.: Schemat modelu Mamba.

Porównując grafiki 2.14 oraz 2.13 można zobaczyć wiele podobieństw. Główne elementy pojawiają się w obu modelach, jednak w przypadku Mamby jest jasno zaznaczony jest mechanizm selekcji, a kolorami jest zaznaczona sekcja pamięci GPU na której są zaalokowane konkretne części modelu.

W przypadku zadań syntetycznych, takich jak kopiowanie, Mamba nie tylko z łatwością rozwiązuje te zadania, ale może również ekstrapolować rozwiązania na bardzo długie sekwencje (ponad milion tokenów), gdzie ich długość znacząco przewyższa tę, na której model został wytrenowany. Mamba osiąga wydajność jakości transformatora w zadaniach modelowania języka. Zapewnia również pięciokrotnie szybszy czas generowania w porównaniu do transformatorów o podobnej wielkości i dorównuje wydajnością transformatorom dwukrotnie większym w testach porównawczych takich jak LAMBDA, Pile lub PIQA.

3. Propozycja rozwiązania

4. Przeprowadzenie eksperymentów

4.1. Prezentacja otrzymanych wyników

Przedstawienie muzyki w pracy stanowi wyzwanie w porównaniu na przykład do wizualnych form sztuki, takich jak obrazy. Muzyka jest medium audytywnym, które oddziałuje na słuch i emocje poprzez dźwięki i rytm, co sprawia, że jej pełne doświadczenie jest trudne do przekazania za pomocą samego tekstu. Opisy muzyczne, notacja i partytury mogą dostarczyć informacji o strukturze i elementach utworu, ale nie oddają one w pełni jego brzmienia i atmosfery. W przeciwieństwie do obrazów, które można bezpośrednio zobaczyć i ocenić na papierze, muzyka wymaga aktywnego odtwarzania, by w pełni ją zrozumieć i docenić. W pracy zostaną przedstawione generacje w dwojaki sposób. W zależności od tego, czy model był wytrenowany na plikach MIDI lub notacji ABC, będzie przedstawiony wynik modelu w formie tekstowej lub *piano roll*. Dodatkowo do każdego z tych przykładów zostanie dodana wygenerowana partytura, jednak należy pamiętać, że w związku z ograniczeniami opisanymi we wcześniejszych częściach pracy, reprezentacja plików MIDI na pięciolinii potrafi nie być idealna.

Na grafice 4.1 przedstawiono wynik modelu [DO OPISANIA JAKIEGO MODELU]. Część oznaczoną kolorem niebieskim należy traktować jako sekwencję wejściową dla modelu, na której podstawie model generuje kolejne tokeny. Dodatkowo należy zwrócić uwagę na podane kody kontrolne w linijce 2, 3 i 4. Reprezentację otrzymanej muzyki zapisanej na pięciolinii przedstawia grafika 4.2.

Zależy zwrócić uwagę, że model został poproszony o wygenerowanie melodii w tonacji C-dur (K:C), jednak jak można zaobserwować na obu reprezentacjach, model zaczął dodawać do dźwięków c krzyżyki ($\sim c$ w notacji ABC). Prawdopodobnie wzięło się to z faktu, że podane jako początek sekwencji dźwięki d, nie są typowym rozpoczęciem melodii w tonacji C-dur. Model oczekiwał podania toniki (pierwszego stopnia skali) ponieważ takie fragmenty są najbardziej popularne i dominowały w zbiorze uczący.

4.2. Sposoby oceny wyników

Celem porównania otrzymanych wyników z dowolnych modeli, jest ocena jakości ich pracy. W kontekście modeli dyskryminujących (ang. *discriminative model*), na przykład klasyfikatorów bądź modeli regresyjnych, ocena ich jakości jest dość prosta, ponieważ istnieje

```

1 X:1
2 S:2
3 B:9
4 E:4
5 L:1/8
6 M:3/4
7 K:C
8 D/2d/2d Ad fa | a/g/e/f/ g/f/e/d/ ^cA | d2 dA df |
9 ag/f/ e^c A2 | D2 d^c de | f2 ef/g/ a/g/f/e/ | d2 df e^c |1 d2 D4 :|
10 2 d2 D2 fg |: a2 g2 f2 | ed ^cd e/f/g | a2 g2 f2 | ed ^cd e/f g |
11 a2 g2 f2 | ed ^cd ef/g/ | a2 g2 fe |1 d2 de fg :|2 d2 D4 ||

```

Rysunek 4.1.: Wynik modelu.



Rysunek 4.2.: Wygenerowana melodia.

predefiniowana, prawdziwa wartość w zbiorze testowym, którą model próbuje przewidzieć. W takim przypadku ewaluacja jakości modelu to porównanie prawdziwych danych, z tymi przewidzianymi przez model. Porównanie odbywa się przez policzenie metryk np. dla modelu klasyfikującego celności, precyzji, $F1-score$ lub dla modelu regresyjnego MSE lub R^2 . W przypadku modeli generacyjnych celem treningu jest jak najlepsza aproksymacja rozkładu prawdopodobieństwa danych $P(X_{data})$ lub w przypadku danych oznaczonych łączny rozkład

$P(X, Y)$. W przypadku dużej wymiarowości danych, obliczenie obiektywnych metryk takich jak *log-likelihood* lub dywergencja Kullbacka-Leiblera (*KLD*) często jest nieobliczalne. Dla człowieka weryfikacja wyników modeli generacyjnych takich jak *text-to-speech* lub *text-to-image* jest trywialnym zadaniem, jednak nie jest to oczywiste zadanie algorytmiczne. Często odwołuje się do subiektywnych metryk takich jak *MOS ang. mean opinion score*, które oblicza się jako średnią opinię np. w skali od 1-5 braną z zazwyczaj niewielkiej grupy ludzi. Niestety metryka ta często nie jest wiarygodna ze względu na jej subiektywności niewielką próbę badawczą. Aby wyeliminować te wady serwisy takie jak *HuggingFace* udostępniają narzędzia dla członków społeczności, które pozwalają na ranking modeli, z nadzieją że zgodnie z prawem wielkich liczb, przy wystarczającej liczbie odpowiedzi, uda się otrzymać w miarę obiektywną ocenę. Przykładowym narzędziem tego rodzaju jest *The TTS Arena*[16], która pozwala na ranking modeli *text-to-speech*.

W przypadku muzyki, istnieje możliwość aby zweryfikować poprawność wygenerowanych sekwencji odwołując się do teorii oraz harmonii muzyki. Istnieje kilka narzędzi takich jak *Chordify*, *Hooktheory* lub *Sibelius*, które pozwalają na analizę harmoniczną utworów, dzięki czemu autorzy mogą w prosty sposób analizować i dobierać progresję danej melodii. Niestety większość takich narzędzi jest płatna i nie pozwala na zautomatyzowaną analizę wielu plików. Dodatkowym problemem jest w analizie harmoniczej, szczególnie prowadzonej przez algorytmy, jest rozróżnienie harmonii wertykalnej oraz horyzontalnej. Rozróżnienie to zostało wythumaczone przez Jacoba Colliera, kilkukrotnego laureata nagród *Grammy*, którego zdaniem akord, który nawet zagrany sam brzmi niepoprawnie, w odpowiednim kontekście i przez odpowiednią progresję w kolejnych fragmentach muzyki, może mieć nadany sens, przez co cała sekwencja nabiera muzycznego piękna[17].

Uwaga 4.1. W mojej pracy prawdopodobnie zostanie zastosowane podejście MOS dla dość niewielkiej grupy ludzi, jednak jeśli trial o oprogramowania pozwoli, spróbuję przynajmniej sprawdzić czy taka analiza pozwala na jakieś sensowne obliczenie metryki

Brak obiektywnych metryk, którymi można się posłużyć podczas ewaluacji modelu jest dodatkowym problemem w momencie *fine-tuning* modelu. Ciężko jest wybrać odpowiednie rozwiązanie architektoniczne lub odpowiednie hiperparametry, kiedy jedyną miarą jakości wygenerowanych danych jest subiektywna opinia programisty lub grupy osób wybranych jako wyrocznia.

4.3. Muzyczny test Turinga

[18]

4.3.1. Testowanie melodii innymi LLM-ami

ChatMusician

ChatMusician to otwarty model językowy, który integruje zdolności muzyczne. Model ten został zaprojektowany na bazie modelu LLaMA2 z 7 miliardów parametrów, jest trenowany i dostrajany z wykorzystaniem notacji muzycznej ABC, traktując muzykę jako drugi język. Dzięki temu *ChatMusician* jest w stanie rozumieć i generować muzykę za pomocą tekstowego tokenizatora, bez konieczności stosowania zewnętrznych struktur neuronowych czy tokenizatorów stworzonych specjalnie dla notacji muzycznej. *ChatMusician* może komponować dobrze zorganizowane, pełnometrażowe utwory muzyczne na podstawie tekstu, akordów, melodii, motywów i form muzycznych, przewyższając wyniki modelu *GPT-4*. Model ten osiąga lepsze wyniki niż *LLaMA2* i *GPT-3.5* w zadaniach związanych z rozumieniem teorii muzyki, co potwierdzono w benchmarku *MusicTheoryBench*, zaprojektowanym do oceny zdolności modeli do rozumienia i rozumowania muzycznego. Model został udostępniony dla użytku publicznego na platformie *HuggingFace*.

Model jest potencjalnym rozwiązaniem problemu z oceną wygenerowanej muzyki. Będąc wytrenowanym na prawdziwych fragmentach muzycznych z całego świata i wielu kultur oraz tekstu zawierającego teorię muzyki, powinien w dość obiektywny sposób oceniać fragmenty melodii. Zapytanie o wcześniej wygenerowany fragment przedstawiono na ??.

Model dość przychylnie spojrzał na ten wygenerowany fragment, jednak jego odpowiedzi są dość ogólne. Nie podejmuje on analizy konkretnych fragmentów i dodatkowo udziela informacji o użytych wielu instrumentach, co nie występuje w przedstawionej mu notacji. W kontekście harmonii używa słów takich jak *well-crafted*, *coherent*, *pleasing*. Słowa te nie są “naukowe” przez co ciężko zinterpretować odpowiedź modelu.

Aby upewnić się nie jest to jest domyślne zachowanie należy poprosić model o ocenę fragmentu, który został specjalnie przygotowany, aby nie brzmiał dobrze dla ludzkiego ucha. Odpowiedź modelu przedstawiono we fragmencie 4.4.

Jak się okazało model również ten fragment opisał bardzo przychylnie, jednak dla żadnego człowieka ta muzyka nie byłaby przyjemna. Niestety model *ChatMusician* nie jest dobrym “nauczycielem” i nie jest w stanie poprawnie ocenić czy dany fragment jest poprawny muzycznie.

GPT-4

Aby przetestować

```

1  PROMPT: I have some music in ABC notation that I generated. Could
2  you analyze it and tell me if it sounds natural? Specifically, I would like
3  your feedback on the overall musicality, the sense of harmony, and whether
4  the melody and chord progressions are coherent and pleasing. Here is the
5  ABC notation of the music:
6  X:1
7  L:1/8
8  M:3/4
9  K:C
10 D/2d/2d Ad fa | a/g/e/f/ g/f/e/d/ ^cA | d2 dA df | ag/f/ e^c A2 |
11 D2 d^c de | f2 ef/g/ a/g/f/e/ | d2 df e^c |1 d2 D4 :|2 d2 D2 fg |:
12 a2 g2 f2 | ed ^cd e/f/g | a2 g2 f2 | ed ^cd e/f g | a2 g2 f2 |
13 ed ^cd ef/g/ | a2 g2 fe |1 d2 de fg :|2 d2 D4 ||
14 ANSWER: Overall Musicality: 1. The piece has a pleasant and flowing
15 rhythm with a steady tempo. It features a mix of different instruments and
16 voices, creating an interesting texture. The melodies are catchy and
17 memorable, making the piece enjoyable to listen to.
18 2. Sense of Harmony: The harmonies in this piece are well-crafted and
19 create a rich and full sound. The use of various chords and progressions
20 adds depth and complexity to the music. The harmonic choices are
21 appropriate for the style and genre of the piece.
22 3. Melody and Chord Progressions: The melody and chord progressions in
23 this piece are coherent and pleasing. They follow a logical progression
24 and fit together well. The melodic lines are well-written and expressive,
25 adding to the emotional impact of the piece. The chord progressions
26 provide a solid foundation for the melodies and add interest and variety
27 to the composition.
28 In conclusion, this piece of music has a strong sense of musicality,
29 harmony, and coherence. The melodies and chord progressions are
30 well-constructed and contribute to the overall appeal of the piece.
31 It is a successful example of its genre and style.

```

Rysunek 4.3.: Pytanie i odpowiedź modelu ChatMusician.

1 **PROMPT:** I have some music in ABC notation that I generated. Could
2 you analyze it and tell me if it sounds natural? Specifically, I would like
3 your feedback on the overall musicality, the sense of harmony, and whether
4 the melody and chord progressions are coherent and pleasing. Here is the
5 ABC notation of the music:
6 X:1
7 M:4/4
8 L:1/8
9 K:C
10 C2 e2 ^g2 _A2 | B3 {g}D1 f2 ^c2 |
11 G2 e1/2 ^g2 _A2 | B3 |
12 C2 e2 ^g2 [_A2C] | B3 D22 |
13 C2 _e2 ^g2 _A2 | =B3 D1 f2 ^c2 ||
14 **ANSWER:** Overall Musicality: The piece has a pleasant and flowing
15 quality to it, with a mix of different rhythms and tempos. The use of
16 dynamics and articulations adds depth and interest to the composition.
17 2. Sense of Harmony: The harmonies used in this piece are mostly consonant
18 and create a sense of stability and balance. There are occasional
19 dissonances that add tension and color to the music. 3. Melody and Chord
20 Progressions: The melodies in this piece are well-crafted and follow a
21 logical progression. They are easy to sing or play and have a clear
22 sense of direction. The chord progressions are also coherent and support
23 the melodies effectively. In conclusion, this piece has a good overall
24 musicality, sense of harmony, and coherence between its melody and
25 chord progressions. It is an enjoyable and engaging composition.

Rysunek 4.4.: Odpowiedź modelu na fałszywy fragment.



Rysunek 4.5.: Specjalnie stworzony “brzydki” fragment muzyki.

5. Zakończenie

1. Podsumowanie.
2. Możliwości dalszego rozwoju.
3. Potencjalne obszary zastosowania pracy.

Bibliografia

- [1] Ashish Vaswani i in. *Attention Is All You Need*. 2023. arXiv: [1706.03762 \[cs.CL\]](#).
- [2] Steve Allen. *Beethoven Symphony No. 7, Movement 2 in ABC*. URL: <https://www.ucolick.org/~sla/abcmusic/sym7mov2.html>.
- [3] Nathan Fradet i in. „MidiTok: A Python package for MIDI file tokenization”. W: *Extended Abstracts for the Late-Breaking Demo Session of the 22nd International Society for Music Information Retrieval Conference*. 2021. URL: <https://archives.ismir.net/ismir2021/latebreaking/000005.pdf>.
- [4] Shangda Wu i in. *TunesFormer: Forming Irish Tunes with Control Codes by Bar Patching*. 2023. arXiv: [2301.02884 \[cs.SD\]](#).
- [5] Nitish Shirish Keskar i in. „CTRL - A Conditional Transformer Language Model for Controllable Generation”. W: *arXiv preprint arXiv:1909.05858* (2019).
- [6] Darrell Conklin. *Bach Chorales*. UCI Machine Learning Repository. DOI: <https://doi.org/10.24432/C5G>
- [7] Curtis Hawthorne i in. „Enabling Factorized Piano Music Modeling and Generation with the MAESTRO Dataset”. W: *International Conference on Learning Representations*. 2019. URL: <https://openreview.net/forum?id=r1lYRjC9F7>.
- [8] Shangda Wu i in. „TunesFormer: Forming Irish Tunes with Control Codes by Bar Patching”. W: *Proceedings of the 2nd Workshop on Human-Centric Music Information Retrieval 2023 co-located with the 24th International Society for Music Information Retrieval Conference (ISMIR 2023), Milan, Italy, November 10, 2023*. Red. Lorenzo Porcaro, Roser Batlle-Roca i Emilia Gómez. T. 3528. CEUR Workshop Proceedings. CEUR-WS.org, 2023. URL: <https://ceur-ws.org/Vol-3528/paper1.pdf>.
- [9] Jacob Devlin i in. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. 2019. arXiv: [1810.04805 \[cs.CL\]](#).
- [10] Alec Radford i in. „Language Models are Unsupervised Multitask Learners”. W: (2019).
- [11] Kaiming He i in. „Deep Residual Learning for Image Recognition”. W: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016, s. 770–778. DOI: [10.1109/CVPR.2016.90](#).

- [12] Albert Gu i in. „Combining Recurrent, Convolutional, and Continuous-time Models with Linear State Space Layers”. W: *Advances in Neural Information Processing Systems*. Red. M. Ranzato i in. T. 34. Curran Associates, Inc., 2021, s. 572–585. URL: https://proceedings.neurips.cc/paper_files/paper/2021/file/05546b0e38ab9175cd9Paper.pdf.
- [13] Albert Gu i in. *HiPPO: Recurrent Memory with Optimal Polynomial Projections*. 2020. arXiv: [2008.07669](https://arxiv.org/abs/2008.07669) [cs.LG].
- [14] Aaron Voelker, Ivana Kajić i Chris Eliasmith. „Legendre Memory Units: Continuous-Time Representation in Recurrent Neural Networks”. W: *Advances in Neural Information Processing Systems*. Red. H. Wallach i in. T. 32. Curran Associates, Inc., 2019. URL: https://proceedings.neurips.cc/paper_files/paper/2019/file/952285b9b7e7a1be5aa7849f32ffff05-Paper.pdf.
- [15] Albert Gu i Tri Dao. *Mamba: Linear-Time Sequence Modeling with Selective State Spaces*. 2023. arXiv: [2312.00752](https://arxiv.org/abs/2312.00752) [cs.LG].
- [16] *The TTS Arena*. URL: <https://huggingface.co/blog/arena-tts>.
- [17] Jacob Collier. *That’s not a wrong note, you just lack confidence*. URL: https://www.youtube.com/watch?v=meha_FCcHbo.
- [18] Christopher Ariza. „The Interrogator as Critic: The Turing Test and the Evaluation of Generative Music Systems”. W: *Computer Music Journal* 33.2 (2009), s. 48–70. ISSN: 01489267, 15315169. URL: <http://www.jstor.org/stable/40301027> (term. wiz. 30.05.2024).