



Akademia Górniczo-Hutnicza im. Stanisława Staszica w Krakowie

Wydział Informatyki

PRACA DYPLOMOWA

Generowanie 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

Streszczenie

Praca poświęcona jest eksploracji możliwości generowania muzyki przy pomocy dużych modeli językowych (*Large Language Models*, LLM). Celem badania jest analiza architektur modeli językowych i ich zastosowanie w procesie tworzenia muzyki. Praca skupia się na podobieństwach strukturalnych między muzyką a tekstem, argumentując, że te dwie formy ekspresji dzielą wspólne cechy sekwencyjne, co pozwala na adaptację modeli językowych do generowania muzycznych kompozycji. W pracy przedstawiono również przegląd istniejących formatów cyfrowej reprezentacji muzyki oraz zbiory danych, które mogą być użyte do treningu modeli generatywnych. Przez przeprowadzenie serii eksperymentów, w których wykorzystano różnorodne architektury, takie jak transformery i modele przestrzeni stanów, przeanalizowano efektywność różnych podejść w kontekście generowania muzyki. Została podjęta próba oceny wygenerowanych fragmentów muzycznych przy pomocy dostępnych modeli LLM, które są udostępniane komercyjnie lub zostały wytrenowane, aby rozumieć muzykę. Praca oferuje nowe perspektywy na temat potencjału LLM w dziedzinie sztucznej kreatywności muzycznej, otwierając drogę do dalszych badań w tej fascynującej przestrzeni interdyscyplinarnej.

Abstract

The work is devoted to explore the possibilities of music generation using large language models (LLM). The aim of the study is to analyze the architectures of language models and their application in the process of music creation. The paper focuses on the structural similarities between music and text, arguing that these two forms of expression share common sequential features, which allows the adaptation of language models for the generation of musical compositions. The paper also provides an overview of existing formats for digital representation of music and datasets that can be used to train generative models. By conducting a series of experiments using a variety of architectures, such as transformers and state space models, the effectiveness of different approaches in the context of music generation was analyzed. An attempt was made to evaluate the generated musical fragments using commercially available LLM models or those trained to understand music. The work offers new perspectives on the potential of LLM in the field of artificial musical creativity, paving the way for further research in this fascinating interdisciplinary space.

Spis treści

1	Wstęp	1
1.1	Motywacja	1
1.2	Cel pracy	2
1.3	Zarys pracy dyplomowej	2
2	Opis aktualnego stanu wiedzy	3
2.1	Podobieństwa pomiędzy muzyką a tekstem	3
2.2	Obecne próby generowania muzyki	4
2.3	Cyfrowa reprezentacja muzyki	5
2.3.1	WAV i MP3	6
2.3.2	Format MIDI	7
2.3.3	Notacja ABC	8
2.3.4	Porównanie reprezentacji utworu muzycznego	11
2.4	Zbiory danych	14
2.4.1	Johann Sebastian Bach Chorales	14
2.4.2	The MAESTRO v3.0	14
2.4.3	IrishMAN	15
2.5	Architektura transformera	15
2.5.1	Sieci RNN	15
2.5.2	LSTM	17
2.5.3	Algorytm uwagi	18
2.5.4	Transformer	20
2.6	Modele przestrzeni stanów	24
2.6.1	Opis modeli <i>state space</i>	24
2.6.2	Mamba	26
3	Przeprowadzenie eksperymentów	29
3.1	Przygotowanie danych	29
3.2	Trening modeli	30
3.3	Prezentacja otrzymanych wyników	34
3.4	Sposoby oceny wyników	38

3.4.1	Testowanie melodii innymi LLM	40
3.4.2	Muzyczny test Turinga	46
4	Zakończenie	49
5	Podziękowania	51
	Bibliografia	53

Spis rysunków

2.1	Fragment chorału J.S. Bacha.	6
2.2	Plik WAV otwarty w programie Audacity.	7
2.3	Muzyka zapisana w pliku MIDI.	8
2.4	Zapis wielu głosów w notacji ABC.	9
2.5	Prosta sekwencja muzyczna.	11
2.6	Reprezentacja muzyki w postaci tokenów.	12
2.7	Przykład utworu pochodzącego ze zbioru MAESTRO.	15
2.8	Przykład wielogłosowego fragmentu ze zbioru IrishMAN.	16
2.9	Schemat budowy fragmentu sieci RNN.	16
2.10	Komórki LSTM połączone między sobą.	18
2.11	Schemat mechanizmu uwagi oraz ich kilkukrotne złożenie.	20
2.12	Schemat transformera.	21
2.13	Algorytm uwagi z maską.	22
2.14	Schemat modelu przestrzeni stanów.	25
2.15	Schemat modelu Mamba.	27
3.1	Czas treningu 100 tysięcy kroków modeli.	32
3.2	Wykresy treningu pierwszych 250 tysięcy kroków modeli.	33
3.3	Czas generowania tokenów przez modele.	34
3.4	Wynik modelu.	35
3.5	Wygenerowana melodia ABC modelem GPT-2 60M	35
3.6	Wielogłosowy fragment melodii wygenerowany modelem Mamba 60M	36
3.7	Wygenerowana melodia MIDI zapisana w notacji muzycznej.	37
3.8	Wygenerowana melodia MIDI modelem GPT2 6M.	37
3.9	Niepoprawny harmonicznie wygenerowany fragment MIDI.	38
3.10	Pytanie dla modelu ChatMusician.	41
3.11	Odpowiedź modelu ChatMusician.	42
3.12	Specjalnie stworzony niepoprawny fragment muzyki.	43
3.13	Odpowiedź modelu ChatMusician na fałszywy fragment.	43
3.14	Odpowiedź modelu GPT-4 na wygenerowany fragment.	44
3.15	Odpowiedź modelu GPT-4 na fałszywy fragment.	45

Spis tabel

2.1	Porównanie różnych zapisów muzyki.	12
3.1	Parametry modelu GPT-2.	31
3.2	Parametry modelu Mamba.	31
3.3	Porównanie wygenerowanych melodii ABC.	39
3.4	Sybiektywna opinia autora na temat modeli MIDI.	40

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 jej tworzenia. Jednakże w erze cyfrowej, której towarzyszy rozwój sztucznej inteligencji, pojawiają się nowe możliwości w dziedzinie generowania 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 generowania 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, aby następnie na podstawie zapytania udzielonego przez użytkownika wygenerować odpowiednią odpowiedź. Modele te znalazły zastosowanie w maszynowej translacji tekstu, analizie sentymentu, odpowiadaniu na pytania użytkownika czy generowania kodu na podstawie poleceń oraz kontekstu nawet dużych projektów. Muzyka swoją strukturą jest bardzo podobna do tekstu. Każda nuta tak jak słowo jest zależna nie tylko od tych, które występują przed nią, ale również od szerszego kontekstu utworu, jak również i tonacji, w której dany utwór został napisany. Dokonując tego porównania, można łatwo zauważyć, dlaczego modele analizujące tekst oraz język pisany są potencjalnie dobrymi kandydatami do próby analizy i generowania muzyki. Obecnie w celach generatywnych stosuje się modele o bardzo dużej ilości parametrów, przez co nie są one często dostępne dla przeciętnego użytkownika, który nie posiada dedykowanego sprzętu. Dodatkowym utrudnieniem jest typowa architektura oparta o transformer, który jest bardzo dobrym wyborem w tego rodzaju zadaniach, jednak generowanie danych przy jego pomocy jest powolne. Obecnie pojawiają się nowe modele do zadań NLP, które jeszcze nie zostały dobrze zbadane w zadaniach muzycznych, a rozwiązują problemy modeli transformerowych. Podejście do generowania muzyki w sposób podobny do metod używanych w generacji tekstu może skutkować powstaniem analogicznych narzędzi, czyli na przykład programów do autouzupełniania melodii. Artyści korzystający

z narzędzi do tworzenia muzyki mogliby otrzymywać podpowiedzi od modelu, który na podstawie już stworzonej przez artystę muzyki, generowałby możliwe jej zakończenia. Osobami, które mogłyby być zainteresowane takimi narzędziami, nie są zawodowi artyści, którzy tworzą muzykę z pasji, a bardziej osoby, które potrzebują na przykład ścieżki dźwiękowej pod gry wideo lub filmy. W ten sposób osoby tylko z hobbistyczną ilością wiedzy otrzymałyby możliwość uczestniczenia w procesie twórczym.

1.2. Cel pracy

Celem pracy jest zbadanie architektur modeli dużych modeli językowych i zastosowania samych modeli w problemie generowania muzyki. Zostanie podjęta analiza możliwych do użycia zbiorów danych i sposobów, w jaki muzyka jest w nich zapisana oraz jakie są wady i zalety każdego z nich. W tym celu zostanie wytrenowane kilka modeli o różnych architekturach i liczbie parametrów na różnym typie danych. Praca poruszy problem weryfikacji wygenerowanych melodii i przedstawi jedno z potencjalnych rozwiązań.

1.3. Zarys pracy dyplomowej

Praca dyplomowa składa się z czterech głównych rozdziałów. W rozdziale pierwszym przedstawiono wstęp do tematu, motywację oraz cel pracy. Rozdział drugi zawiera szczegółowy opis aktualnego stanu wiedzy na temat generowania muzyki przy pomocy dużych modeli językowych. Omówiono w nim podobieństwa pomiędzy muzyką a tekstem, przegląd istniejących prób generowania muzyki, cyfrowe reprezentacje muzyki oraz różne zbiory danych używane w eksperymentach. Szczególną uwagę poświęcono architekturze transformera oraz modelom przestrzeni stanów, które stanowią podstawę dla omawianych modeli generatywnych. W rozdziale trzecim opisano metodologię przeprowadzenia eksperymentów, w tym przygotowanie danych, proces treningu modeli oraz zaprezentowano otrzymane wyniki. Przeanalizowano także różne sposoby oceny wygenerowanych melodii, w tym testowanie ich innymi dużymi modelami językowymi. Ostatni rozdział zawiera podsumowanie wyników pracy, wnioski oraz sugestie dotyczące dalszych badań w dziedzinie generowania muzyki przy użyciu sztucznej inteligencji.

2. Opis aktualnego stanu wiedzy

2.1. Podobieństwa pomiędzy muzyką a tekstem

Tekst i muzyka wykazują fundamentalne podobieństwo strukturalne, opierające się na wzorach i sekwencjach w celu przekazywania i wywoływania emocji. Oba wykorzystują hierarchiczną organizację: zdania tworzą akapity, podczas gdy nuty i akordy tworzą melodie i harmonie. Podobnie jak tekst jest zbudowany zgodnie z prawami składni i gramatyki, muzyka przestrzega zasad rytmu i harmonii. Muzyka posiada poza wymiarem czasowym (poziomym) również aspekt pionowy, czyli możliwość wystąpienia kilku nut, które są zagrane w tym samym czasie w formie akordu. Jest to element, który nie występuje w tekście zapisanym alfabetem łacińskim, przez co istotne jest, aby forma przetwarzania notacji muzycznej uwzględniała możliwość ujęcia tego zjawiska. Każda nuta w akordzie spełnia swoje harmoniczne zadanie, co powoduje, że nie każda grupa nut zagrana razem brzmi przyjemnie dla ludzkiego ucha. Melodie składają się z wielu akordów następujących po sobie zgodnie z regułami harmonii. Poza nutami w zapisie muzycznym są również pauzy, znaki artykulacji i akcentów, jak i zapis rytmu każdej nuty. Analizując opisane porównania, użycie narzędzi przetwarzania języka naturalnego (NLP) jest potencjalnie dobrym wyborem w celu analizy muzyki.

Duże modele językowe (LLM) wykazały w ostatnich latach niezwykle możliwości w zadaniach przetwarzania języka naturalnego i całkowicie zrewolucjonizowały to pole [1, 2]. Ich sukces doprowadził do znacznego rozwoju badań w tym kierunku. LLM, będące biegłe w rozpoznawaniu i generowaniu sekwencji tekstowych, można dostosować do generowania muzyki poprzez uczenie ich w analogiczny sposób struktur w kompozycji muzycznej. Trenując na obszernych zbiorach danych partytur muzycznych, LLM mogą przewidywać i tworzyć spójne sekwencje nut, akordów i rytmów, tworząc w ten sposób całkowicie nowe utwory muzyczne, które zachowują integralność stylistyczną i strukturalną.

Przyczynowe modelowanie języka *Causal Language Modeling* (CLM), znane również jako auto-regresyjne modelowanie języka, jest techniką wykorzystywaną w przetwarzaniu języka naturalnego do przewidywania następnego słowa w sekwencji na podstawie poprzednich słów. Model generuje tekst krok po kroku, gdzie przewidywanie bieżącego słowa jest uwarunkowane wcześniej wygenerowanymi słowami. Oznacza to, że model wykorzystuje własne wcześniejsze wyniki jako dane wejściowe do przewidywania następnego słowa. Na każdym etapie model generuje rozkład prawdopodobieństwa nad każdym unikalnym tokenem, wskazując,

jak prawdopodobne jest, że dany token będzie następny w sekwencji. Typowo w tego rodzaju modelach występuje pojęcie okna kontekstowego, czyli ilość wstecznych słów, które model bierze pod uwagę w celu generowania kolejnego.

$$\begin{aligned} &\text{sekwencja wejściowa: } x_1, x_2, x_3, \text{ gdzie } x_i \in X \\ &P_{t_1}(X) = \text{Model}(x_1, x_2, x_3) \\ &x_4 \sim P_{t_1}(X) \\ &P_{t_2}(X) = \text{Model}(x_2, x_3, x_4) \\ &x_5 \sim P_{t_2}(X) \end{aligned} \tag{2.1}$$

Przykładowe podejście do auto-regresyjnego generowania sekwencji pokazano w równaniu 2.1. Zapis P_{t_n} oznacza rozkład prawdopodobieństwa nad zbiorem możliwych elementów sekwencji X w kroku czasowym t_n , z którego jest następnie losowany kolejny komponent sekwencji. Algorytm jest agnostyczny w stosunku do typu danych, więc w zbiorze X mogą być zarówno wszystkie słowa, litery, jak i elementy notacji muzycznej. Model z równania działa na oknie czasowym o długości trzech tokenów, jednak jest to parametr, który można zmienić w czasie tworzenia modelu. Mechanizm generowania może działać w nieskończoność, jednak zazwyczaj kończy on się w momencie wygenerowania sekwencji o konkretnej długości ustalonej przez użytkownika lub kiedy zostanie napotkany element zbioru, który oznacza koniec sekwencji. W przypadku tekstu może być to kropka, natomiast w muzyce podwójna kreska taktowa.

2.2. Obecne próby generowania muzyki

Jedną z pierwszych prób generowania muzyki przy pomocy algorytmów datuje się na rok 1957. Powstały utwór o nazwie *Illiac Suite* został napisany na kwartet smyczkowy przy użyciu metod Monte Carlo oraz procesów Markowa [3]. Na użycie sieci neuronowych oraz uczenia maszynowego, jakie znamy dzisiaj, należało poczekać do roku 2014. Organizacją, która podjęła się tego tematu jest jeszcze obecnie działająca Magenta [4]. Projekt ten skupia się na badaniu i rozwijaniu algorytmów, które mogą wspierać artystów w tworzeniu nowych dzieł poprzez generowanie muzyki za pomocą technik uczenia maszynowego. Magenta stworzyła różne modele generatywne, które potrafią komponować nowe utwory muzyczne. Jednym z modeli jest *MusicVAE*, który bazując na sieciach konwolucyjnych, był w stanie generować sekwencje muzyczne, oraz interpolować pomiędzy nimi [5]. Innym modelem stworzonym przez tę organizację jest *GanSynth* [6], który był w stanie generować naturalnie brzmiącą muzykę przy pomocy modelu *GAN* (*generative adversarial network*). Przełomowym momentem w rozwoju modeli generatywnych było wprowadzenie modeli transformera w roku 2017 [7]. Doprowadziło to do wzmożonego zainteresowania generatywnym tworzeniem muzyki. Rok później pojawił się model *Music Transformer*, który implementował architekturę *decoder-*

only i był trenowany na zbiorze danych MAESTRO [8]. Był on w stanie generować długie sekwencje muzyczne o długotrwałej strukturze, jednak melodia generowana przez narzędzie była dość monotonna, kiedy istniało niedostateczne dopasowanie do zbioru treningowego lub niewystarczająco zróżnicowane dane szkoleniowe. Ograniczało to przydatność narzędzia do złożonych zadań generowania muzyki i wymagało dodatkowej optymalizacji lub szkolenia z większymi i bardziej zróżnicowanymi zestawami danych w celu poprawy jakości generowanej muzyki [9]. Bardzo podobnym modelem jest *Museformer* [10], który implementując inny rodzaj algorytmu *attention* w dekodzie rozwiązuje problemy swojego poprzednika z rozmiarem sekwencji. Idea polega na tym, że nie trzeba skupiać się na całej sekwencji muzycznej z tym samym poziomem ważności. W modelu zostały wyróżnione dwa rodzaje algorytmu uwagi, czyli *fine-grained attention* (drobnoziarnista) oraz *coarse-grained attention* (gruboziarnista). Wynika to z obserwacji mówiącej, że dla generowanej muzyki ważniejsze są bliższe sobie tokeny niż te, które znajdują się w większej odległości. W zależności od tego, na podstawie jakich tokenów dokonujemy obliczeń, używa się innego algorytmu.

Obecnie najczęściej wykorzystywanymi modelami do generowania muzyki są modele *text-to-music*. Modele te działają na zasadzie udzielenia zapytania do modelu, który na jego podstawie generuje odpowiedni fragment. Typowe zapytanie może brzmieć na przykład „*Slow tempo, bass-and-drums-led reggae song. Sustained electric guitar. High-pitched bongos with ringing tones. Vocals are relaxed with a laid-back feel, very expressive*”. Przykładem takiego modelu jest *MusicLM* stworzony przez Microsoft [11]. Zwrócił on na siebie uwagę dzięki zdolności do tworzenia wysokiej jakości muzyki, która może trwać nawet kilka minut i charakteryzuje się wysoką jakością dźwięku. Niestety model nie został udostępniony jako *open source*, natomiast zbiór danych zawierający ponad pięć tysięcy przykładów uczących wraz z ich szczegółowym opisem można pobrać ze strony *Kaggle*. Istnieją modele stworzone przez społeczność, które powstały w oparciu o zaproponowaną przez Microsoft architekturę i są dostępne do własnego użytku. Jedną z głównych zalet *MusicLM* jest to, że może generować złożoną muzykę, która strukturą i spójnością jest podobna do muzyki tworzonej przez ludzkich kompozytorów. Oznacza to, że model ma potencjałbycia wykorzystanym w różnych zastosowaniach, od generowania tła muzycznego do filmów i gier po wspomaganie kompozytorów w procesie twórczym.

2.3. Cyfrowa reprezentacja muzyki

W kulturze zachodniej nuty zapisuje się na pięcioliniach ułożonych jedna pod drugą. Jeśli muzyka jest wielogłosowa, pierwsze takty z każdej pięciolinii są łączone przy pomocy linii lub okalane nawiasem klamrowym. Na pięcioliniach umieszczone są takie elementy jak nuty, pauzy, oznaczenia dynamiki, tempa i inne znaki dotyczące zapisu, w jaki sposób należy wykonać dany utwór. Przykładowy zapis muzyki znajduje się na ilustracji 2.1. Przedstawiony fragment składa się z czterech głosów brzmiących naraz oraz z czterech taktów. Utwór pocho-

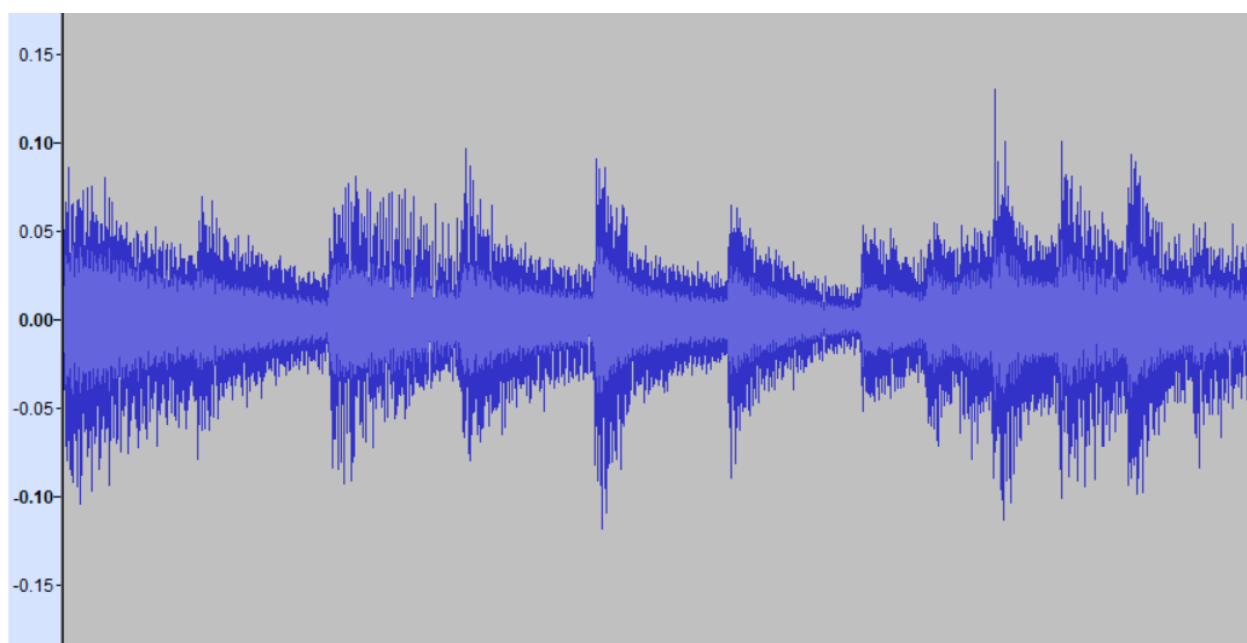


Rysunek 2.1.: Fragment chorału J.S. Bacha.

dzi ze zbioru chorałów Jana Sebastiana Bacha, o którym to zbiorze jest mowa w dalszej części pracy. W związku z tym, że zbiór zawiera pliki MIDI, dokonano automatycznej translacji fragmentu do zapisu nutowego przy pomocy narzędzia MuseScore, przez co zapis nutowy zawiera dość nietypowe elementy, takie jak zmianę klucza w środku taktu. Na początku muzykę zapisywano na papierze, jednak po pojawieniu się komputerów zapis nutowy, jak i nagranie utworu 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 zapisu muzyki w postaci plików był *MusicXML*, który przy pomocy zwykłych znaczników XML reprezentował notację muzyczną.

2.3.1. WAV i MP3

WAV (*ang. Waveform Audio Format*) to format plików binarnych, znany z możliwości zapisywania dźwięku bez użycia jakiegokolwiek algorytmu kompresji [12]. Dzięki temu pliki WAV charakteryzują się najwyższą jakością dźwięku, jednak ich rozmiary są bardzo duże, co może być problematyczne przy przechowywaniu dużych zbiorów danych na dysku. Plik WAV jest najwierniejszą cyfrową reprezentacją dźwięku analogowego. Jego duży rozmiar wynika z faktu, że dźwięk jest zapisywany z częstotliwością 44,1 kHz, czyli w każdej sekundzie w pamięci zapisywane jest 44100 próbek. Format ten został stworzony w 1991 roku i od tego czasu stał się jednym z najbardziej rozpowszechnionych formatów audio, obsługiwanych przez praktycznie każde oprogramowanie do edycji dźwięku. W celu wizualizacji dokonano elektronicznego odtworzenia utworu przedstawionego w postaci nutowej na ilustracji 2.1,



Rysunek 2.2.: Plik WAV otwarty w programie Audacity.

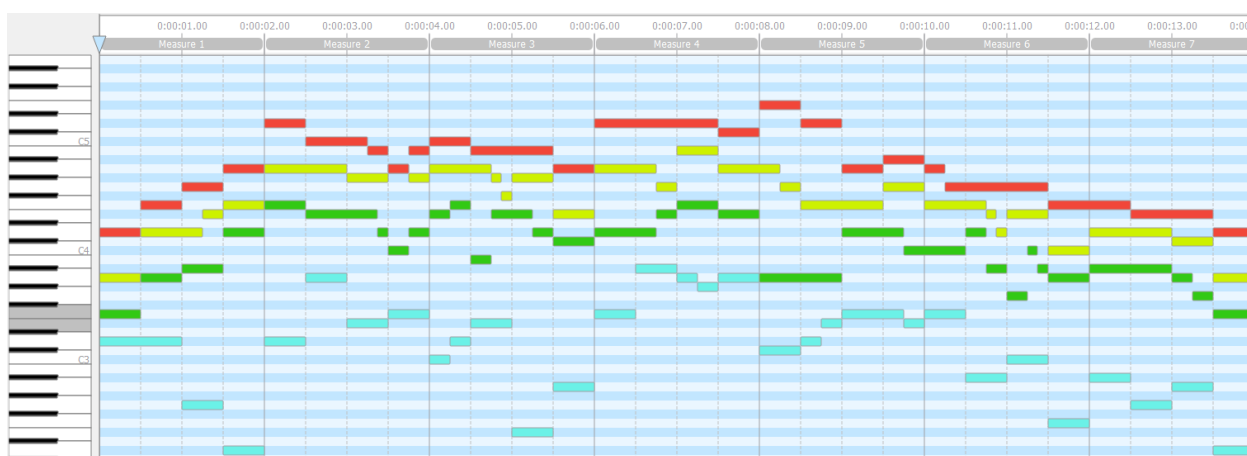
co zaowocowało powstaniem fali dźwiękowej, przedstawionej na grafice 2.2.

Często nie jest konieczne przechowywanie dźwięku w formacie bezstratnym, ponieważ większość informacji można zachować przy użyciu mniejszej ilości danych. Jednym z najpopularniejszych standardów kodowania stratnego jest format MP3. Kompresja zmniejsza dokładność kodowania i usuwa częstotliwości, które nie są słyszalne dla człowieka. Stratne kodowanie stara się znaleźć równowagę między jakością dźwięku a rozmiarem pliku. W kontekście uczenia maszynowego zmniejszona ilość próbek przy zachowaniu większości informacji jest zjawiskiem pożądanym. Dzięki temu przynajmniej częściowo rozwiązujemy problem związany z „przekleństwem wymiarowości”, czyli trudnością w przetwarzaniu danych o bardzo wysokiej liczbie cech (wymiarów). Stratne kodowanie pozwala na efektywniejsze zarządzanie danymi, co może być kluczowe w procesach analitycznych i modelowaniu.

2.3.2. Format MIDI

MIDI (*ang. Musical Instrument Digital Interface*) to standard opisujący protokół komunikacji, interfejs cyfrowy oraz złącze, które umożliwiają połączenie elektronicznych instrumentów, komputerów oraz innych urządzeń muzycznych [13]. Pojedynczy kabel MIDI może przesyłać informacje na temat szesnastu kanałów jednocześnie, z których każdy może reprezentować inny instrument. Każda interakcja z instrumentem, taka jak naciśnięcie klawisza

czy szarpnięcie struny, jest zapisywana jako „zdarzenie” (*event*), które zawiera takie dane jak znacznik czasu, wysokość dźwięku oraz jego głośność. Dane z urządzeń MIDI są zapisywane w specjalnych plikach z rozszerzeniem `.mid` lub `.midi`. Umożliwiają one przechowywanie, rozpowszechnianie oraz edytowanie dźwięków w specjalnie stworzonym do tego oprogramowaniu. Ponieważ plik MIDI nie przechowuje rzeczywistej fali dźwiękowej nagrałej przez mikrofon, możliwa jest np. późniejsza zmiana instrumentu, który będzie odtwarzał zapisane dźwięki. Typowym sposobem przedstawienia danych MIDI jest tak zwany *piano roll*, który można porównać do dwuwymiarowego układu współrzędnych, gdzie na osi poziomej reprezentowany jest czas, a oś pionowa przedstawia konkretne dźwięki, pokazane jako klawisze fortepianu. Przykładowe porównanie zapisu MIDI z zapisem nutowym przedstawiono na grafikach 2.3 oraz 2.1.



Rysunek 2.3.: Muzyka zapisana w pliku MIDI.

Warto zauważyć, że zapis pliku w formacie MIDI, w porównaniu do zapisu nutowego, umożliwia zdecydowanie większą ekspresję, ponieważ jest to zapis odtworzenia przez artystę pewnego utworu. Pozwala to na ujęcie zdecydowanie więcej emocji, jak i osobowości samego muzyka, co może zostać zaobserwowane na przykład poprzez lekkie przeciąganie pewnych nut lub zmianę głośności ich odegrania. Informacje na temat tego, w jaki sposób odegrać dany utwór (akcenty, głośność nut), mogą zostać zapisane w tradycyjnej notacji, jednak nawet wtedy artysta może lekko modyfikować zamysł kompozytora.

2.3.3. Notacja ABC

Notacja ABC to system zapisywania nut w formie tekstowej przy użyciu znaków ASCII. Notacja ta pojawiła się w latach 70. XX wieku w celu zapisu i nauki tradycyjnych irlandzkich melodii [14]. W kolejnej dekadzie system ten został rozwinięty przez Chrisa Walshaw,

który używał go do zapisywania tradycyjnych melodii, zanim opanował standardowy zachodni zapis nutowy. Walshaw stworzył program *abc2mtex*, który na podstawie notacji ABC generował komendy umożliwiające zapis partytur w formacie *MusicTex*. Obecnie używanym standardem notacji ABC jest wersja z 2011 roku.

```
X:1
Q:1/4=120
V:1
L:1/16
M:4/4
K:C clef=G2
D4F4G4A4|d4c6B2A2B2|c4B8A4|d12^c4|
V:2
L:1/16
M:4/4
K:C clef=G2
A,4D6E2F4|A8^G4A2^G2|A6^G^F^G4E4|A6G2B4A4|
V:3
F,4A,4^A,4D4|F4E7DC2D2|E2F2B,2E4D2^C4|D6E2F4E4|
V:4
L:1/16
M:4/4
K:C clef=F4
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.

Rysunek 2.4 przedstawia fragment muzyki 2.1 zapisany w notacji ABC. Format w dość zwięzły sposób zapisuje partyturę. Poza konkretnymi nutami i rytmem w tym formacie można zapisać również dodatkowe informacje na temat utworu. W każdej linii zaczynającej się od znaku A-Z, a następnie dwukropka, znajduje się tak zwane „pole informacyjne”. W tych polach można zapisywać takie informacje, jak tytuł utworu, metrum oraz wiele innych danych, w tym z jakiego zbioru muzycznego pochodzi utwór lub jaki jest jego tytuł. Wiele z tych informacji powinna zostać usunięta w procesie wczesnego przetwarzania danych, aby nie zagłuszały tekstem najistotniejszych pól w utworze. Do ważnych pól należą przede wszystkim: domyślna długość nuty (L:), metrum (M:) oraz tonacja w jakiej utwór został napisany (K:). Chociaż model jest w stanie „odgadnąć” te informacje na podstawie dźwięków, ich wyraźne podanie daje modelowi więcej informacji na temat każdej z sekwencji. Dodatkową zaletą podawania tych informacji podczas treningu jest możliwość użycia ich jako

początkowej sekwencji, na podstawie której model będzie dalej generował melodię. Dzięki temu można mieć kontrolę nad takimi aspektami, jak tonacja i metrum utworu, co pozwala na bardziej precyzyjne kierowanie procesem generowania muzyki.

Notacja ABC wspiera również melodie polifoniczne przy pomocy znaczników V:. Ich ilość nie jest ograniczona, w związku z czym 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 partii instrumentów [15].

Tokenizacja plików MIDI

Tokenizacja to proces przekształcania danych na mniejsze jednostki zwane tokenami, którymi w kontekście tekstu mogą być słowa, ich fragmenty lub nawet pojedyncze znaki. Tokenizacja jest jednym z pierwszych kroków w analizie tekstu, umożliwiającym dalsze przetwarzanie, takie jak analizę składniową, semantyczną czy ekstrakcję cech na podstawie zbioru treningowego. Proces ten pomaga w standaryzacji, eliminując niejednoznaczności i ułatwiając dalsze kroki przetwarzania, takie jak stemming czy lematyzację. Algorytmy NLP operujące na tokenach są zazwyczaj bardziej efektywne i łatwiejsze do implementacji. Przetwarzanie mniejszych jednostek tekstu zmniejsza złożoność obliczeniową i pozwala na bardziej precyzyjną analizę [16]. Jednym z głównych powodów tokenizacji jest potrzeba konwersji danych tekstowych na reprezentację numeryczną, która może być przetwarzana przez algorytmy uczenia maszynowego. Korzystając z tej reprezentacji, można trenować modele do wykonywania różnych zadań, takich jak klasyfikacja, analiza sentymentu lub generowanie sekwencji tekstowych.

W kontekście plików MIDI proces tokenizacji jest niezbędny, ponieważ format MIDI zawiera skomplikowane struktury i informacje binarne, które są trudne do bezpośredniego wykorzystania przez model. Celem tego procesu jest utworzenie sekwencji tokenów, które będą reprezentować wartości atrybutów nut (wysokość, wartość, czas trwania) lub zdarzenia czasowe zawarte w melodii.

Token może przyjmować jedną z trzech form:

- nazwa tokenu – słowna reprezentacja zdarzenia MIDI, np. *Pitch_50*;
- id – unikalna wartość liczbowa przypisana konkretnemu zdarzeniu;
- bajt – unikalna wartość liczbowa przypisana podczas treningu tokenizera.

Nowe słownictwo jest zapisywane w postaci *look up table* łączącej nazwę tokenu z odpowiadającym jej id lub bajtem. Trening tokenizera polega na obliczeniu kodowania gramatycznego, 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 partii danych 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, dokonano analizy popularnego algorytmu *REMI*. **REvamped MIDI** reprezentuje zdarzenia jako sekwencje tokenów wysokości tonu, dynamiki (głośności), długości oraz czasu trwania poprzez kolejne tokeny. 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. Każda nuta jest reprezentowana przez sekwencję tokenów:

- Pos – relatywną pozycję nuty w stosunku do początku taktu zapisaną jako liczba z przedziału 0-31;
- Tempo – w jakim tempie (*bpm*) nuta została zagrana;
- Pitch – wysokość dźwięku;
- Vel – głośność zapisana jako liczba z przedziału 0-127;
- Dur – długość dźwięku zapisana jako ułamek ćwierćnuty.

Dodatkowymi tokenami widocznymi na ilustracji są tokeny **Bar** odpowiadające za początek taktu oraz **Rest** kodujące pauzy lub przerwę w muzyce. Wiele algorytmów tokenizacji jest zaimplementowanych w bibliotece MidiTok dla języka Python [17].



Rysunek 2.5.: Prosta sekwencja muzyczna.

2.3.4. Porównanie reprezentacji utworu muzycznego

W tabeli 2.1 zostało przedstawione porównanie zapisu całego utworu, którego fragment przedstawiono na grafice 2.1. Kolumna „długość sekwencji” oznacza liczbę wartości numerycznych opisujących całą sekwencję w danym kodowaniu. W plikach WAV i MP3 są to



Rysunek 2.6.: Reprezentacja muzyki w postaci tokenów.

	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.

amplitudy sygnału dźwiękowego, w MIDI – długość sekwencji tokenów, a w ABC – liczba znaków tekstowych.

Jak widać, zapis pliku w postaci fali dźwiękowej zajmuje najwięcej miejsca na dysku oraz długość sekwencji, która reprezentuje cały utwór jest największa. W celu uzyskania sekwencji z pliku MIDI został użyty tokenizer REMI, o którym jest mowa w rozdziale 2.3.3. Podczas używania innych tokenizatorów długość sekwencji może różnić się nieznacznie pomiędzy nimi. Różnica w rozmiarze na dysku oraz długością sekwencji pomiędzy zapisem MIDI oraz ABC nie jest duża, jednak warto zwrócić uwagę, że gdyby w oryginalnym utworze pojawiły się powtórzenia sekcji melodii, notacja ABC byłaby jeszcze bardziej oszczędna od plików MIDI. Powodem tego jest fakt, że notacja ABC zapisuje bezpośrednio zapis nutowy w postaci tekstowej, którego elementami są znaki repetycji. Są zapisywane znakami :| oraz |: , a w pliku MIDI cała powtarzana sekwencja musi zostać ponownie odegrana.

Do głównych zalet plików MIDI należy ich wszechstronność. Nie są one ograniczone do sztywno określonego rytmu w postaci ćwierćnut, ósemek czy szesnastek, jak w innych notacjach muzycznych, co daje zdecydowanie większe możliwości ekspresji. Dodatkowo, przy użyciu narzędzi komputerowych, bardzo łatwo można zamienić dowolny zapis nutowy (np. MusicXML lub ABC) na plik MIDI. Niestety, w drugą stronę, czyli konwersja pliku MIDI na zapis nutowy, jest znacznie trudniejsza, głównie z powodu nieokreślonego rytmu. Istnieją narzędzia, takie jak MuseScore, które potrafią wygenerować partyturę z pliku MIDI, jednak skuteczność tych narzędzi nie zawsze jest wysoka, zwłaszcza w przypadku bardziej skomplikowanych utworów, gdzie rytm i artykulacja mogą być trudne do dokładnego oszacowania. Mimo tych wyzwań, pliki MIDI pozostają niezwykle użyteczne w zastosowaniach muzycznych, od tworzenia i edytowania muzyki po jej analizę i odtwarzanie.

Zapis pliku w postaci fali dźwiękowej niesie ze sobą wiele wad. Jednak jest to często jedyne

rozwiązanie, aby zebrać dużą ilość danych, szczególnie w przypadku muzyki nowoczesnej. Wówczas twórcy nie udostępniają zazwyczaj plików MIDI ani zapisu nutowego, na podstawie którego powstał dany utwór. Z tego powodu praca z plikami MP3 lub czasochłonna translacja muzyki na inny format jest często konieczna.

Zaletą plików dźwiękowych w zapisie cyfrowym jest dostęp do potencjalnie istniejącego zapisu odśpiewanego tekstu. Pomimo że notacja ABC pozwala na zapis tekstu w znaczniku W:, jest to jedynie zapis słowny, a nie faktyczne odśpiewanie. W związku z tym pliki WAV są bardziej odpowiednie do próby replikacji czyjegoś głosu, ponieważ zawierają rzeczywisty zapis audio, który można analizować i przetwarzać.

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 mniejszą 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 procesowania danych. Podejście takie zostało zaproponowane w modelu *Tunesformer* [18], który wziął przykład z artykułu [19], w którym autorzy dodawali do korpusu treningowego danych tekstowych takie znaczniki, jak *Books*, *Horror*, *Relationships*, *Legal*. Oznaczenia te były podawane na początku pytania podawanego dla modelu, przez co model miał w jasny sposób podane, w jakim stylu powinien udzielić odpowiedzi. W przypadku muzyki w zapisie ABC zostały wprowadzone nowe kody kontrolne:

- S: – reprezentuje liczbę sekcji w całym utworze, znacznik łatwy do policzenia ponieważ początek i koniec taktu jest jawnie oznaczany przy pomocy symboli [|, ||, |], |: , :: i :|;
- B: – oznacza liczbę taktów w danej sekcji, zlicza tylko wystąpienia znaku |;
- E: – 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 według 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 sprawia, że notacja staje się nieprawidłowa. Tak

samo jak w przypadku znaczników, które opisują elementy zapisu, jak tonacja czy metrum. Służą one nam oraz modelowi jako wskazówki, którymi powinien się sugerować podczas generowania muzyki.

2.4. Zbiory danych

2.4.1. Johann Sebastian Bach Chorales

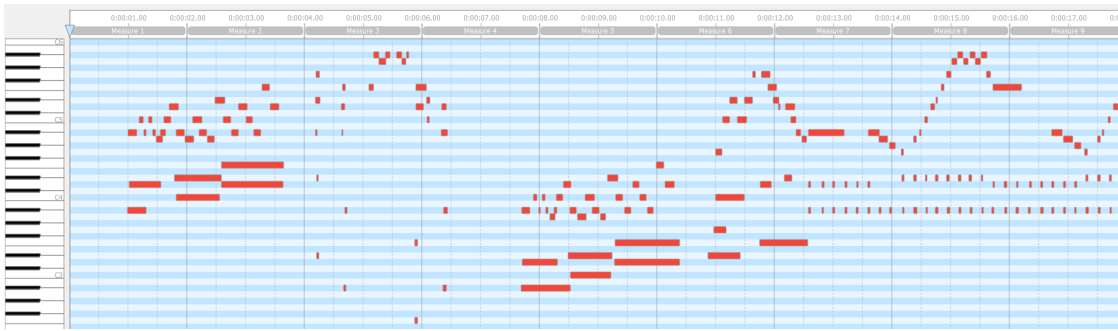
JSB Chorales [20] to zbiór krótkich, czterogłosowych utworów pierwotnie skomponowanych przez Jana Sebastiana Bacha w XVIII wieku. Napisał je, biorąc wcześniej istniejące melodie ze współczesnych mu hymnów luterkańskich, a następnie harmonizując je w celu stworzenia pozostałych trzech głosów. Zbiór danych zawiera 382 takich chorałów. Dodatkowo ich rytm został znormalizowany, przez co najmniejszą długością nuty jest szesnastka. Dzięki temu uproszczeniu zbiór jest bardzo łatwo konwertowalny do dowolnego formatu 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 są wysokości kolejno brzmiących dźwięków. 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 [21], jednak w tym zbiorze ten problem nie występuje.

2.4.2. The MAESTRO v3.0

Zbiór danych MAESTRO [22] powstał w współpracy z *Minnesota International Piano-e-Competition*, czyli międzynarodowym konkursem pianistycznym. Litera *e* w nazwie odnosi się do fortepianów Yamaha Disklavier używanych podczas konkursu, które poza tradycyjną funkcjonalnością są obudowane elektronicznymi sensorami, które pozwalają między innymi zapis odegranej muzyki w formacie MIDI. Z tego powodu organizacja Magenta działająca wewnątrz Google użyła zapisu melodii z wielu edycji konkursu, aby stworzyć zbiór danych zawierający około 200 godzin muzyki fortepianowej. W tych 200 godzinach znajduje się 1276 występów pianistycznych, w skład których wchodzi ponad 7 milionów nut. Dwa rodzaje zbioru, jakie można pobrać, to zbiór w wersji WAV oraz MIDI. Pierwszy z nich, z powodów omawianych we wcześniejszych częściach pracy, zajmuje aż 122 GB, natomiast wersja MIDI zajmuje tylko 81 MB.

Porównując zapis prostego utworu widocznego na grafice 2.3 z tym na grafice 2.7, osoby nawet nie znające się na muzyce są w stanie wywnioskować, że jest to utwór bardziej złożony. Jak widać, odegrane nuty często nie padają na sztywno określone rozpoczęcie taktu,



Rysunek 2.7.: Przykład utworu pochodzącego ze zbioru MAESTRO.

przez co potencjalnie wygenerowana muzyka na podstawie tego zbioru będzie wydawała się bardziej ekspresyjna oraz „ludzka”.

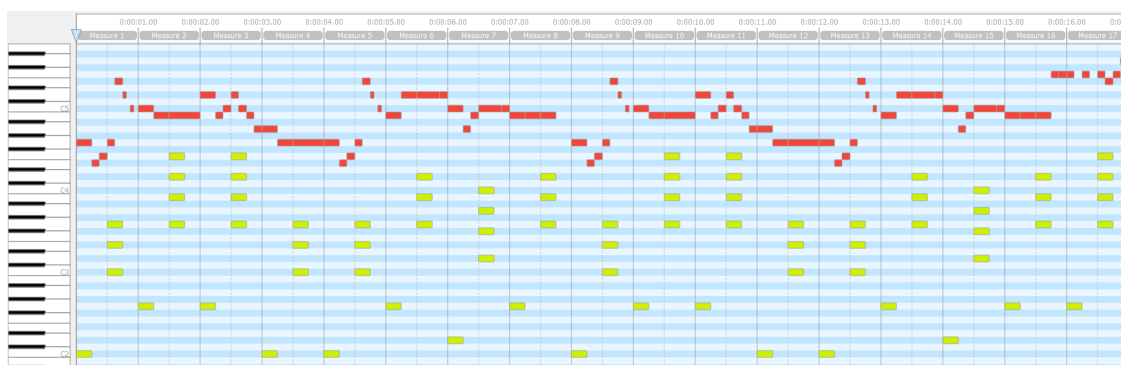
2.4.3. IrishMAN

Zbiór danych IrishMAN (*ang. Irish Massive ABC Notation*) [23] zawiera 216,248 fragmentów tradycyjnych irlandzkich melodii. Pochodzą one ze stron takich jak thesession.org oraz abcnotation.org, które znane są z udostępniania tradycyjnej muzyki w różnych formatach. Aby zapewnić jednolitość formatowania, najpierw wszystkie utwory zostały przekonwertowane do notacji MusicXML, a następnie do notacji ABC. Istnieją bliźniacze zbiory IrishMAN-MIDI oraz IrishMAN-XML, które zawierają te same utwory, jednak zapisane w odpowiadających ich nazwom formatach. Wszystkie fragmenty muzyki należą do domeny publicznej, w związku z czym można korzystać z tych melodii bez obaw o łamanie praw autorskich. 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 generatywnych.

2.5. Architektura transformera

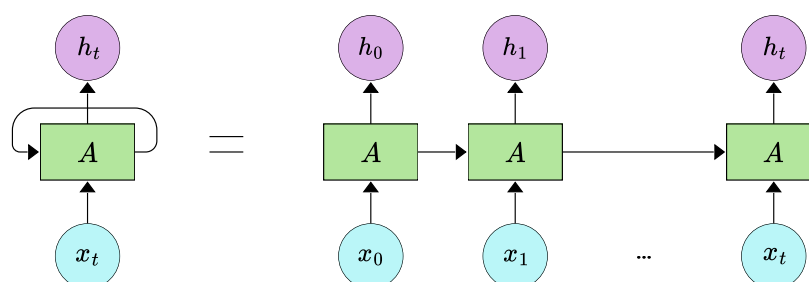
2.5.1. Sieci RNN

Pierwszymi próbami użycia sieci neuronowych w przetwarzaniu sekwencji było użycie rekurencyjnych sieci neuronowych (RNN). Architektura ich była prosta i dodawała ona pętlę



Rysunek 2.8.: Przykład wielogłosowego fragmentu ze zbioru IrishMAN.

w taki sposób, aby wartości wynikowe mogły być przekazywane pomiędzy krokami. Rysunek 2.9 przedstawia fragment sieci oraz jej rozwinięcie.



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 w zdaniu „Rok temu ukończyłem kurs lotnika ... Teraz staram się o licencję *pilota*” wymagane jest szersze spojrzenie na całą strukturę wypowiedzi i „przypomnienie” informacji z początku zdania. W przypadku muzyki sytuacja wygląda dokładnie tak samo. Przykładowo, aby rozwiązać akord dysonansowy na konsonansowy lub dominantowy na tonikę, wymagana jest jedynie informacja na temat ich dźwięków, tak aby dobrać odpowiedni nowy akord. Jednak, jeśli w muzyce przewija się pewien temat, 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 jest to bardzo rzadko osiągalne.

Dodatkowym problemem podczas treningu modeli jest problem zanikającego gradientu. Równanie 2.3 przedstawia obliczenia przeprowadzane w celu wyznaczenia konkretnego stanu ukrytego h_i .

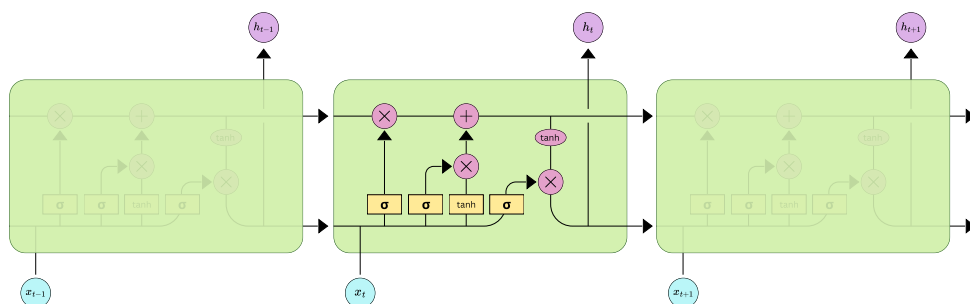
$$\begin{aligned}h_{i+1} &= Wh_i + Zx_i \\h_1 &= Wh_0 + Zx_1 \\h_2 &= W^2h_0 + WZx_1 + Zx_2 \\&\vdots \\h_N &= W^Nh_0 + W^{N-1}Zx_1 + \dots + Zx_N\end{aligned}\tag{2.3}$$

W powyższym równaniu wartość N oznacza długość okna kontekstowego, na którym operuje model, W – macierz wag z obecnego kroku, a Z jest wagami z kroku poprzedniego. Aby obliczyć stan ukryty, na podstawie którego zostanie dokonana dalsza predykcja, należy policzyć każdy poprzedni stan zaczynając od pierwszego. Aby uchwycić odpowiedni kontekst fragmentu testu, długość okna kontekstowego powinna być maksymalnie długa. Z tego powodu podczas obliczeń macierz wag W będzie podnoszona do bardzo wysokiej potęgi, przez co jej wartości poniżej 1 będą zanikać do 0 (zanikający gradient), a te będące powyżej 1 będą dążyły do nieskończoności (eksplodujący gradient). W przypadku drugiego problemu istnieją metody regularyzacyjne takie jak „ucinienie gradientu” (*gradient clipping*), które starają się go rozwiązać [24]. W przypadku zanikającego gradientu można zastosować normalizację pakietu danych (*batch normalization*), która jednak nie zawsze rozwiązuje ten problem [25].

2.5.2. LSTM

Aby rozwiązać omawiane problemy sieci rekurencyjnych została zaproponowana architektura sieci nazwana *Long Short Term Memory* (LSTM) [26]. Ich budowa pozwala na długofalowe zapamiętywanie i przekazywanie informacji. Warstwy LSTM są zazwyczaj łączone między sobą, co pozytywnie wpływa na działanie sieci. Komórka LSTM zwraca dwa wektory. Pierwszy, stan ukryty (*hidden state*), jest wektorem przechowującym informacje o bieżącym stanie sieci w danym kroku. Jego wartości są aktualizowane przy każdym kroku czasowym na podstawie nowych danych wejściowych oraz poprzedniego stanu ukrytego. Reprezentuje on krótkoterminowe informacje w sekwencji, które są potrzebne do generowania odpowiedzi na bieżącym etapie. Stan komórki (*cell state*) to dodatkowy wektor, który przechowuje informacje długofalowe, czyli takie, które nie wynikają bezpośrednio z niedalekich poprzednich komórek. Jego struktura umożliwia łatwe przekazywanie informacji w dłuższej perspektywie, co jest kluczowe w zadaniach wymagających pamięci długoterminowej. Stan komórki jest modyfikowany za pomocą trzech bramek: bramki wejścia (*input gate*), bramki zapomnienia

(*forget gate*) oraz bramki wyjścia (*output gate*) [27].



Rysunek 2.10.: Komórki LSTM połączone między sobą.

Główną częścią budowy komórki jest górna część reprezentowana przez prostą, biegnącą przez całą jej długość strzałkę widoczną na ilustracji 2.10, która posiadając jedynie proste operacje liniowe, przekazuje stany ukryte z jednej części do drugiej. Dolna część komórki zawiera wcześniej wspomniane bramki, które kontrolują, które informacje są ważne i mają być zachowane, a które można zignorować. Są one zbudowane z warstw sigmoidalnych oznaczonych na grafice literami σ . Wartościami zwrotnymi bramek jest liczba pomiędzy 0 a 1, gdzie mała wartość znaczy niską, a wysoka – dużą wagę informacji.

Architektura takiego modelu sprawia, że jest on o wiele mniej podatny na problem znikającego gradientu, jednak jest to sytuacja, która niestety może nadal występować [28]. 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*). Algorytm generuje rozwinięcie sieci rekurencyjnej na wiele kroków czasowych, tworząc rozłożony na warstwy graf obliczeń. Następnie stosowana jest tradycyjna propagacja wsteczna do obliczenia gradientów błędów w stosunku do wag sieci na każdym kroku czasowym. Gradienty te są sumowane, a wagi są aktualizowane w celu minimalizacji błędów predykcji w całej sekwencji, umożliwiając sieci lepsze dopasowanie do danych sekwencyjnych.

Istnieje kilka wariantów [sieci LSTM](#), np. dwukierunkowe LSTM, które pozwalają na przepływ informacji w dwóch kierunkach na raz, a nie tylko z lewej strony do prawej [29]. 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.

2.5.3. Algorytm uwagi

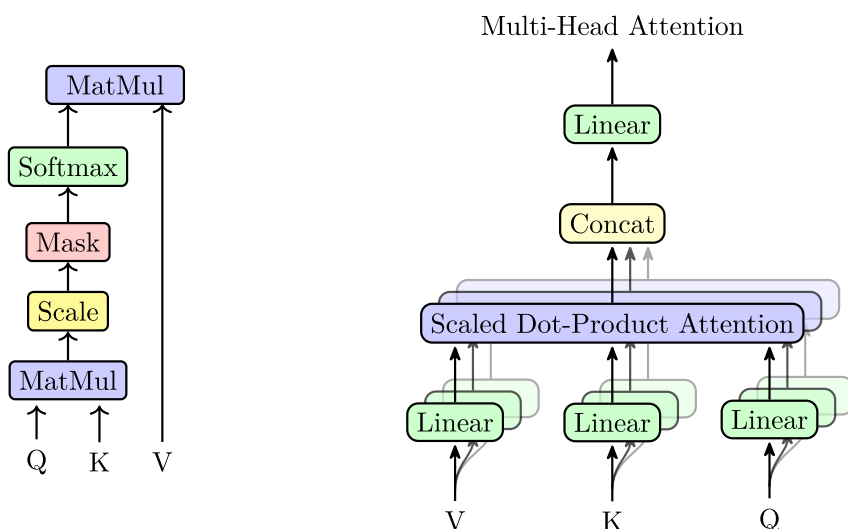
Mechanizm uwagi (*ang. attention*) został zaproponowany w przełomowym artykule „*Attention is All You Need*” [7]. Algorytm ten działa w ustalonym oknie kontekstowym (*ang.*

context window), w którym przydziela każdemu słowu wagę reprezentującą jego ważność w danym kontekście. Metoda ta imituje sposób w jaki człowiek, czytając na przykład książkę, skupia się na różnych fragmentach tekstu w danym momencie, interpretując dane słowa w kontekście innych, już przeczytanych oraz tych, które zostaną przeczytane za chwilę. Mechanizm uwagi został opracowany, aby rozwiązać problem występujący w rekurencyjnych sieciach neuronowych, gdzie „ważność” słowa była tym większa, im bliżej końca sekwencji się ono znajdowało. Sprawiało to, że modele często ignorowały informacje znajdujące się na początku okna kontekstowego. Mechanizm uwagi został zaprojektowany tak, aby identyfikować najwyższe korelacje między słowami w zdaniu, zakładając, że zostaną one wykryte na podstawie zbioru treningowego. Korelacja ta jest przechwytywana w wagach neuronów poprzez propagację wsteczną albo z samonadzorowanego szkolenia wstępnego, albo z nadzorowanego dostrajania. Algorytm działa na wartościach liczbowych, więc przed analizą każda sekwencja musi zostać zakodowana (*embedded*) w postaci wektora o określonej długości.

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

Podstawowa wersja algorytmu nazywana jest *self-attention* lub *scaled dot-product attention*, ponieważ sekwencja „zwraca uwagę” na samą siebie, czyli innymi słowy każdy token w sekwencji jest rozważany z każdym innym poza sobą. 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 dla danego słowa są wszystkie inne. Jest to zdecydowanie lepsze podejście niż to z RNN lub LSTM, gdzie zależności 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.

Przedstawione na rysunku 2.11 wartości K (Key), V (Value) i Q (Query) w mechanizmie *self-attention* są uzyskiwane przez przemnożenie danych wejściowych i odpowiadających im nauczonych macierzy wag. Dla każdego elementu sekwencji (reprezentowanego przez jego zapytanie Q) obliczane jest podobieństwo do każdego innego elementu sekwencji (reprezentowanego przez klucze K). Najczęściej używa się iloczynu skalarnego, który jest normalizowany przez wymiar osadzonego wektora (d_k), aby uzyskać współczynniki uwagi, a następnie, aby wartości sumowały się do jedynki, używana jest funkcja softmax. Ten krok można rozumieć jako tworzenie słownika, w którym kluczami są poszczególne tokeny, a odpowiadającymi im wartościami są miary ważności w całej sekwencji. Każda wartość V jest ważona przez współczynniki uwagi obliczone w poprzednim kroku, co pozwala modelowi skupić się na najistotniejszych elementach sekwencji. Cały ten złożony proces sprowadza się do jednego równania zapisanego w 2.4. W algorytmie można dodać opcjonalną maskę, która umożliwia



Rysunek 2.11.: Schemat mechanizmu uwagi oraz ich kilkukrotne złożenie.

zasłonięcie niektórych wartości uwagi, dzięki czemu te tokeny nie będą na siebie wpływać. Prosty schemat tego algorytmu został pokazany na lewej stronie grafiki 2.11.

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V \quad (2.4)$$

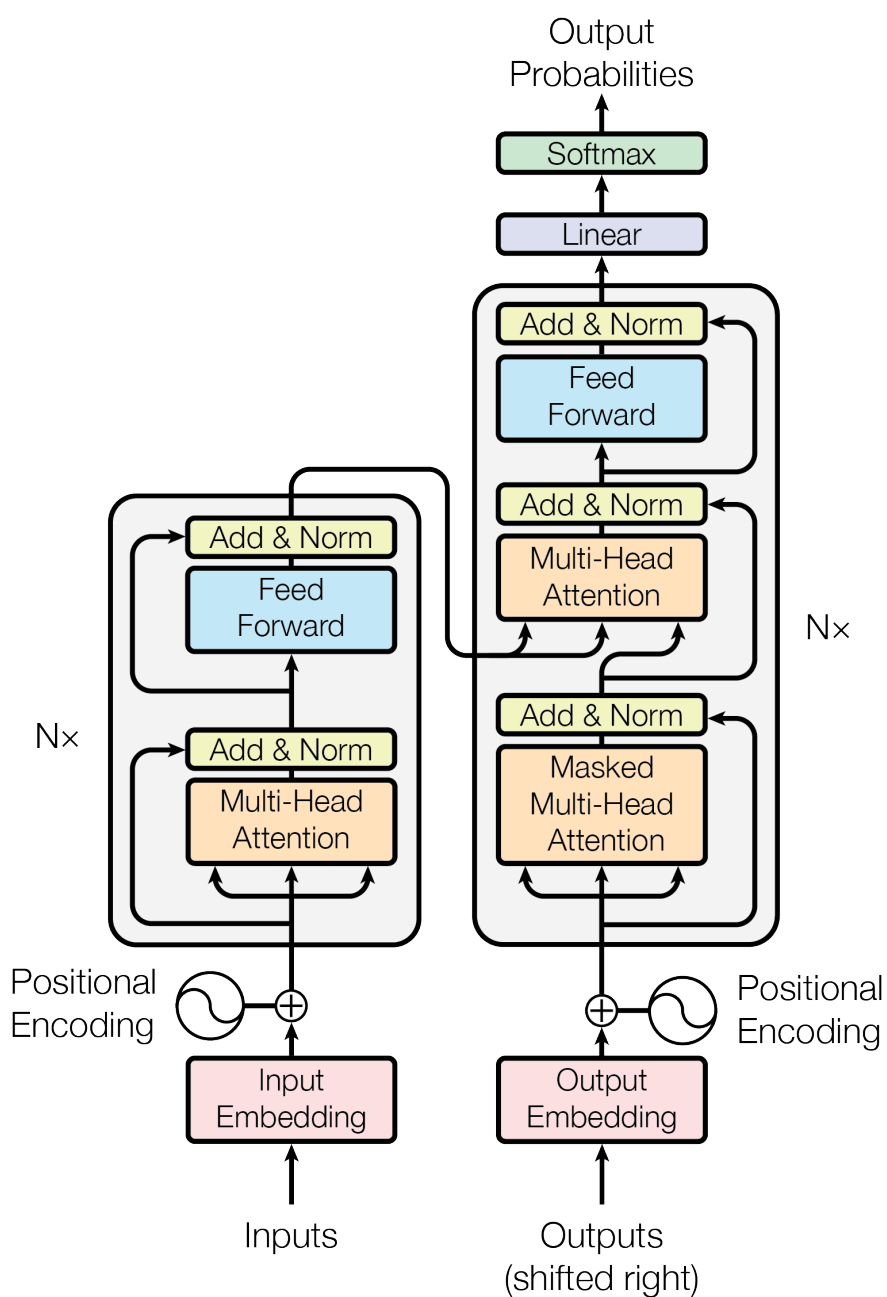
Aby zeskalować model i zwiększyć liczbę jego parametrów w celu użycia go do bardziej zaawansowanych celów, można użyć algorytmu *multi-headed attention*. Jest to kilka warstw uwagi, przed wejściem do których znajdują się warstwy liniowe. Wartości po wyjściu z algorytmu *attention* są konkatelowane, a następnie ponownie przetwarzane przez warstwę liniową. Tym sposobem, zamiast tylko jednej instancji mechanizmu uwagi, model może korzystać z wielu. Schemat takiego podejścia pokazano na prawej stronie grafiki 2.11.

2.5.4. Transformer

Wraz z wprowadzeniem algorytmu uwagi, w tym samym artykule został zaproponowana architektura modelu uczenia maszynowego wykorzystująca *multi-headed attention*. Cały model składa się z dwóch części nazwanych enkoderem oraz dekoderelem.

Enkoder

Zadaniem enkodera jest zapoznanie się z całym korpusem tekstu i przekazanie zebranych informacji do dekodera. Mechanizm samo-uwagi *self-attention* jest centralnym elementem



Rysunek 2.12.: Schemat transformera.

każdej warstwy enkodera. Pozwala on modelowi analizować relacje między różnymi tokenami w sekwencji niezależnie od ich odległości od siebie. Warto zwrócić uwagę, że enkoder

jest dwukierunkowy, czyli tokeny są w stanie „patrzeć w przód” analizowanej sekwencji. Drugi komponent każdego enkodera to sieć neuronowa typu *feed forward*, która składa się z dwóch w pełni połączonych warstw z funkcjami aktywacji *ReLU*. Każda z dwóch głównych części składowych jest otoczona mechanizmem dodania rezydualnego, poprawiającego wydajność treningu oraz warstwą normalizacji. Dane przed wejściem są kodowane przy pomocy tabeli kodowań do wektorów, a następnie, aby zapewnić modelowi informacje na temat kolejności słów we fragmencie, jest dodawany wektor kodujący ich kolejności. Omawiane w sekcji 2.5.3 pracy wartości K (klucze) oraz V (wartości) otrzymane na wyjściu enkodera zostają przekazane do dekodera.

Dekoder

Zadaniem dekodera 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 górną z wartościami $-\infty$, którą następnie dodaje się do wartości uwagi co zostało pokazane na grafice 2.13. Maska zawiera takie wartości z tego powodu, że jedną z operacji jaką wykonuje funkcja softmax jest e^{x_i} , co w przypadku wartości $x_i = -\infty$ da zero, co sprawi, że te wartości zostaną zignorowane w dalszych operacjach modelu.

	Ala	ma	kota	i	psa		Ala	ma	kota	i	psa
Ala	0.13	0.18	0.16	0.15	0.18		Ala	0.13	$-\infty$	$-\infty$	$-\infty$
ma	0.68	0.02	0.08	0.14	0.02		ma	0.68	0.02	$-\infty$	$-\infty$
kota	0.06	0.25	0.14	0.11	0.23	→	kota	0.06	0.25	0.14	$-\infty$
i	0.21	0.14	0.16	0.17	0.14		i	0.21	0.14	0.16	0.17
psa	0.27	0.11	0.16	0.18	0.12		psa	0.27	0.11	0.16	0.18

Rysunek 2.13.: Algorytm uwagi z maską.

Drugi mechanizm uwagi umożliwia dekodерowi skupienie się na odpowiednich częściach

sekwencji wejściowej (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 wyjścia enkodera. Następnie, tak jak w poprzednim przypadku, występuje warstwa *feed forward* i warstwy normalizujące. W ostatniej części modelu 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.

Istnieją modele transformerowe niekorzystające z całej architektury. Przykładem *encoder-only* transformera jest model BERT [30], który jest typowo wykorzystywany do zadań, np. klasyfikacji tekstu. W przypadku takiego zadania nie potrzeba omawianej wcześniej maski, ponieważ model powinien analizować kontekst całego fragmentu w „obie strony”. Wystarczy dołożyć na wyjście modelu warstwę liniową o odpowiednim rozmiarze oraz zastosować odpowiednią funkcję aktywacji, aby przerobić enkoder na model klasyfikacyjny. Analogicznie, jeśli celem zadania jest wyłącznie generowanie tekstu lub muzyki, enkoder nie jest konieczny. W takim problemie należy usunąć drugą warstwę uwagi (bez maski), ponieważ wektory kluczy i wartości pochodzące z enkodera będą nieznane. Generacja kolejnych tokenów odbywać się będzie jedynie na podstawie podanej sekwencji lub tokenu BOS (*beginning of sequence*), co da modelowi największą swobodę 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 enkoderami z dodaną maską w pierwszej warstwie *attention*. Jest to często wykorzystywana „sztuczka”, którą można wykorzystać w wielu bibliotekach implementujących gotowe części składowe modelu transformera.

Trening transformera może być **w znaczny** sposób zrównoleglony, co wynika z faktu, że mechanizm uwagi jest w stanie niezależnie liczyć wartości uwagi pomiędzy słowami. Przykładowo w zdaniu „Ala ma kota.” wartości pomiędzy *ma* i *kota* można policzyć, nie znając tych pomiędzy *Ala* i *ma*. Niestety, podczas generowania kolejnych tokenów na podstawie istniejącej już sekwencji należy policzyć ponownie wartości uwagi, przez co dla sekwencji o długości n należy wykonać n^2 operacji. Jest to 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, ponieważ stan ukryty zawiera informacje o **wszystkich** poprzednich stanach, jednak w praktyce nigdy nie ma dostępu do takich informacji.

2.6. Modele przestrzeni stanów

2.6.1. Opis modeli *state space*

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

Typowo dla każdej jednostki czasu t model:

1. dokonuje kodowania sekwencji wejściowej $x(t)$;
2. oblicza reprezentację ukrytą $h(t)$;
3. przewiduje sekwencję wynikową $y(t)$.

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

$$\text{równanie stanu: } h'(t) = Ah(t) + Bx(t) \quad (2.5)$$

$$\text{równanie wynikowe: } y(t) = Ch(t) + Dx(t) \quad (2.6)$$

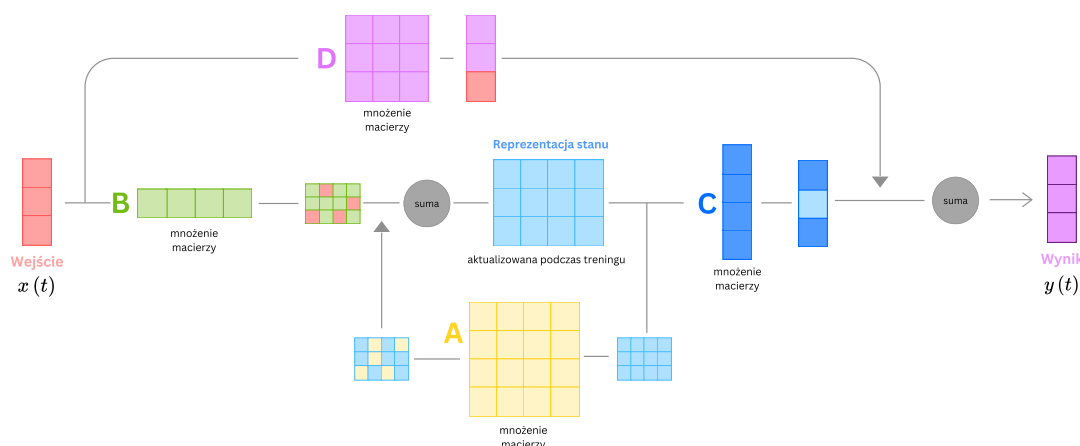
Celem jest poznanie reprezentacji stanu $h(t)$ w taki sposób, aby przejść z sekwencji wejściowej do sekwencji wynikowej.

Równanie stanu 2.5 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 2.6 jest opisanie, jak obecny stan ukryty wpływa na sekwencję wynikową oraz jak dane wejściowe wpływają na wynik. Te informacje są zapisywane odpowiednio w macierzy C oraz D . Schematyczne przedstawienie działania modelu zostało przedstawione na grafice 2.14 [31].

Macierz D przypomina połączenia rezydualne lub połączenia typu *skip-connection* w głębokich sieciach neuronowych, które zapewniają inną ścieżkę dla danych, aby dotrzeć do ostatnich części sieci neuronowej poprzez pominięcie niektórych warstw [32]. Dla uproszczenia zapisu często pomija się je w zapisie modelu.

Dyskretyzacja modelu

Modele tego typu działają domyślnie w domenie danych ciągłych. Tekst lub muzyka są przykładami danych dyskretnych, przez co model musi zostać do nich 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



Rysunek 2.14.: Schemat modelu przestrzeni stanów.

jest dodatkowym parametrem nazywanym krokiem. W związku z tym, że model na wyjściu zwraca sygnał ciągły, aby zamienić go na dyskretny, należy próbkować go z odpowiednią długością kroku. Po dyskretyzacji modelu można nazwać go modelem *sequence to sequence*. Aby lepiej zobrazować dokonane zmiany, należy spojrzeć na równania 2.7 oraz 2.8.

$$\text{równanie stanu: } h_k = Ah_{k-1} + Bx_k \quad (2.7)$$

$$\text{równanie wynikowe: } y_k = Ch_k \quad (2.8)$$

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 z zapisu 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. Modele przestrzeni stanów 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 trójwymiarowe, należy zastosować jądro o jednym wymiarze. Jednowymiarowa konwolucja działa poprzez przesuwanie filtra o ustalonej szerokości wzdłuż jednowymiarowego sygnału wejściowego, obliczając iloczyn skalarny wartości filtra i fragmentów sygnału w każdej pozycji. Wynik tej operacji tworzy nowy sygnał wyjściowy, który zawiera

informacje o lokalnych cechach oryginalnego sygnału. W tym przypadku jądrem konwolucji (funkcji splotu) są odpowiednio przemnożone macierze A , B oraz C , które są aktualizowane podczas treningu. Tego rodzaju interpretacja modelu jest o tyle wygodna, że pozwala zrównoleglić trening modelu, a co za tym idzie, zmniejszyć jego czas [33]. Kiedy model jest już wytrenowany i gotowy do generowania, 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 wykorzystuje jedynie ostatni stan ukryty. Aby zapewnić dobrą kompresję kontekstu jako A , stosuje się macierz HiPPO [34]. Zadaniem macierzy *High-order Polynomial Projection Operators* jest jak najlepsze skompresowanie wejścia do wektora cech. Konstrukcję macierzy pokazano w równaniu 2.9. W tym przypadku można założyć, że macierz jest kwadratowa.

$$A_{nk} = - \begin{cases} (2n+1)^{1/2}(2k+1)^{1/2} & \text{dla } n > k \\ n+1 & \text{dla } n = k \\ 0 & \text{dla } n < k \end{cases} \quad (2.9)$$

Zbudowanie macierzy w taki sposób prowadzi do zdecydowanie lepszych wyników niż zainicjalizowanie jej losowo. Macierz HiPPO reprodukuje sekwencje w taki sposób, że dane, które znajdują się 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.

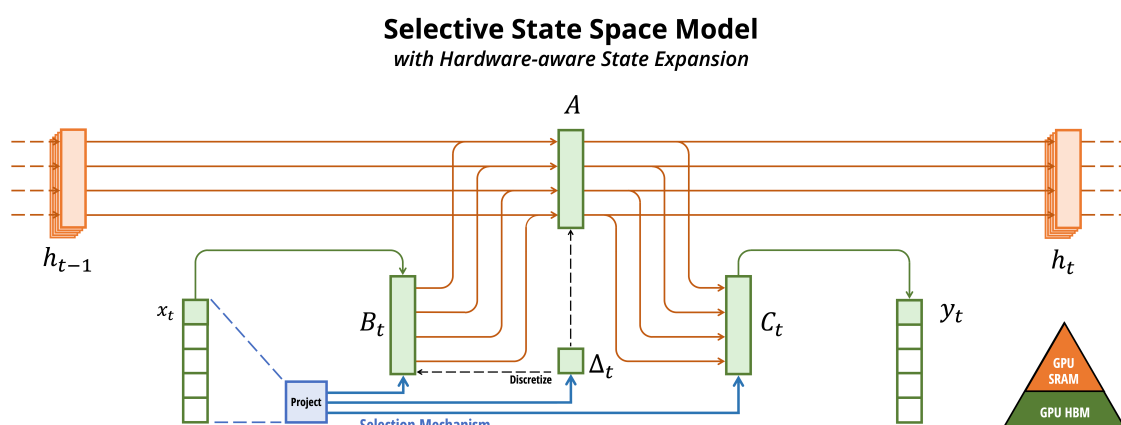
Modele 4S

Po dodaniu omawianych usprawnień do tradycyjnego modelu SSM (dyskretyzacja, macierz HiPPO) nowo powstały model otrzymał 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 i konwolucyjnej rozwiązuje problemy między innymi złożoności obliczeniowej treningu oraz wykorzystywania modelu podczas predykcji.

2.6.2. Mamba

Model Mamba posiada architekturę, która bazuje na modelach 4S, jednak wprowadza [opisane poniżej usprawnienia](#).

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 operowania modelu.



Rysunek 2.15.: Schemat modelu Mamba.

Porównując grafiki 2.15 oraz 2.14 można zobaczyć wiele podobieństw. Główne elementy pojawiają się w obu modelach, jednak w przypadku Mamby jest jasno zaznaczony mechanizm selekcji oraz dyskretyzacja modelu oznaczona jako Δ_t . Model Mamba jest zoptymalizowany pod kątem specyficznych możliwości kart graficznych (GPU), na których jest on typowo uruchamiany i trenowany. Sekcje pamięci, na których są zaalokowane konkretne części modelu, na grafice są zaznaczone kolorem pomarańczowym i zielonym.

W przypadku zadań syntetycznych, takich jak kopiowanie sekwencji, Mamba nie tylko z łatwością rozwiązuje te zadania, ale może również ekstrapolować to rozwiązanie na bardzo długie sekwencje (ponad milion tokenów), których długość znacząco przewyższa te, na której model został wytrenowany. Mamba osiąga wydajność jakości transformera w zadaniach modelowania języka, zapewniając pięciokrotnie krótszy czas generowania w porównaniu do modeli transformerowych o podobnej wielkości i dorównuje wydajnością transformatorom dwukrotnie większym w testach porównawczych [35].

3. Przeprowadzenie eksperymentów

W celu oceny skuteczności modeli generowania muzyki przeprowadzono szereg eksperymentów, w których porównano modele transformerowe oraz Mamba. Eksperymenty te były przeprowadzone na plikach muzycznych w formatach MIDI i ABC, a modele były trenowane w dwóch wariantach wielkościowych: z około 60 milionami parametrów oraz z 6 milionami parametrów. Celem tych badań było zrozumienie, jak różne podejścia i rozmiary modeli wpływają na jakość generowanej muzyki.

3.1. Przygotowanie danych

Pliki w formacie ABC, dzięki swojej tekstowej naturze, nie wymagały skomplikowanego przetwarzania wstępnego. Format ABC zapisuje nuty i inne informacje muzyczne w postaci tekstu, co idealnie pasuje do modeli przystosowanych do pracy z sekwencjami tekstowymi.

Do tokenizacji plików MIDI zastosowano algorytm REMI, o którym była mowa w rozdziale 2.3.3. Przekształca on dane muzyczne w sekwencje tokenów, które mogą być interpretowane przez modele.

Aby jeszcze bardziej zredukować złożoność i objętość danych, zastosowano technikę *Byte Pair Encoding* (BPE). BPE jest algorytmem kompresji tekstu, który jest szeroko stosowany w przetwarzaniu języka naturalnego do redukcji liczby unikalnych tokenów w sekwencjach tekstowych [36]. W kontekście przetwarzania muzyki BPE pomaga w optymalizacji danych wejściowych, co z kolei prowadzi do skrócenia sekwencji wejściowej modelu i przekłada się na lepszej jakości predykcje [37].

Algorytm można opisać w następujących krokach:

1. Inicjalizacja – każdy unikalny token (w kontekście muzyki będzie to nuta, pauza, zmiana tempa itp.) jest traktowany jako pojedynczy symbol;
2. Zliczanie par – algorytm identyfikuje najczęściej występujące pary sąsiadujących tokenów w sekwencji;
3. Łączenie par – najczęściej występujące pary są łączone w jeden nowy symbol, redukując tym samym liczbę symboli w sekwencji;
4. Powtarzanie procesu – kroki 2 i 3 są powtarzane aż do osiągnięcia zdefiniowanego rozmiaru słownika (*vocab size*).

Aby zilustrować działanie algorytmu BPE, można rozważyć prostą sekwencję liter `abcbcab`, na której zostanie przeprowadzony ten algorytm. Na początku należy zliczyć ilość par oraz zidentyfikować te najczęściej powtarzające się. W tym przypadku jest to para `ab`. Zostanie ona zamieniona na nowy token, np. `X`. Nowa sekwencja `XcXcX` zostaje poddana ponownej analizie, z czego w tym przypadku wynika, że najczęstszą parą jest `Xc`. Zostanie ona zamieniona na nowy token `Y`, z czego powstanie sekwencja `YYX`. W przypadku tekstu, tokenami, na podstawie których obliczane są pary, są kolejne bajty tekstu. W kontekście ztokenizowanego zbioru danych MIDI pary będą obliczane na podstawie tokenów, takich jak: `Pitch_C3`, `Velocity_70`, `Duration_1.0`. Jeśli tego rodzaju sekwencja będzie występowała wystarczająco często, przy użyciu BPE zostanie ona zastąpiona jednym tokenem.

W przeprowadzonych eksperymentach ustawiono rozmiar słownika na 20,000, co zapewniło odpowiednią równowagę między kompresją a zachowaniem istotnych informacji muzycznych.

3.2. Trening modeli

HuggingFace to wiodąca firma w dziedzinie przetwarzania języka naturalnego i uczenia maszynowego, znana z opracowywania narzędzi i bibliotek ułatwiających tworzenie, wdrażanie i udostępnianie modeli uczenia maszynowego. Jednym z jej flagowych produktów jest bardzo popularna biblioteka typu *open-source* Transformers, która zapewnia łatwy dostęp do szerokiej gamy wstępnie wytrenowanych modeli dla różnych zadań NLP, takich jak klasyfikacja tekstu, tłumaczenie, podsumowywanie i odpowiadanie na pytania. Biblioteka obsługuje modele takie jak BERT, GPT, T5 i wiele innych, oferując zarówno podstawową architekturę, jak i wstępnie wytrenowane wagi, co znacznie skraca czas i zasoby obliczeniowe potrzebne do dostrajania modeli. HuggingFace udostępnia wysokopoziomowe API pozwalające na tworzenie oraz trening modeli przy pomocy klas, takich jak `Trainer`, `AutoModelForCausalLM` czy `ModelConfig`.

Model reprezentujący transformer został oparty o architekturę modelu GPT-2 [38], a jego konfiguracja została dostosowana, aby posiadał odpowiednią liczbę parametrów. Oryginalny model został stworzony przez OpenAI i jest to obecnie ich najnowszy model, którego wagi, jak i architektura są udostępniane jako *open source*. Jest oparty o architekturę *decoder-only*. Pomimo tego, że wytrenowane modele na danych tekstowych są gotowe do pobrania z platformy HuggingFace, w celu przeprowadzenia rzetelnych badań, modele były trenowane od podstaw wyłącznie na danych muzycznych.

Celem treningu modeli typu *Causal Language Model* (*CausalLM*) jest jak najwierniejsze przewidzenie następnego tokena w sekwencji, biorąc pod uwagę poprzednie tokeny. Podczas szkolenia tokeny wejściowe są wprowadzane do modelu, a model przewiduje rozkład prawdopodobieństwa nad wektorem tokenów, a następnie określa, z takim prawdopodobieństwem losowany jest kolejny token. Funkcja straty (*loss*) jest obliczana na podstawie przewidywań modelu i rzeczywistych tokenów docelowych, które są tokenami wejściowymi przesuniętymi

o jedną pozycję. Typowo używaną funkcją w tego rodzaju problemach jest funkcja entropii krzyżowej *cross-entropy loss*, która mierzy różnicę między dwoma rozkładami prawdopodobieństwa: prawdziwym rozkładem (rozkładem docelowym) a rozkładem przewidywanym przez model.

W poniższych tabelach zostały pokazane konfiguracje modeli, jak i liczba parametrów trenowanych modeli.

Tabela 3.1.: Parametry modelu GPT-2.

Parametr	model 6M	model 60M
n_embed	128	512
n_layer	2	6
n_head	1	4
n_inner	256	2048

Tabela 3.2.: Parametry modelu Mamba.

Parametr	model 6M	model 60M
hidden_size	128	512
state_size	4	8
num_hidden_layers	8	12

Parametry modelu GPT-2:

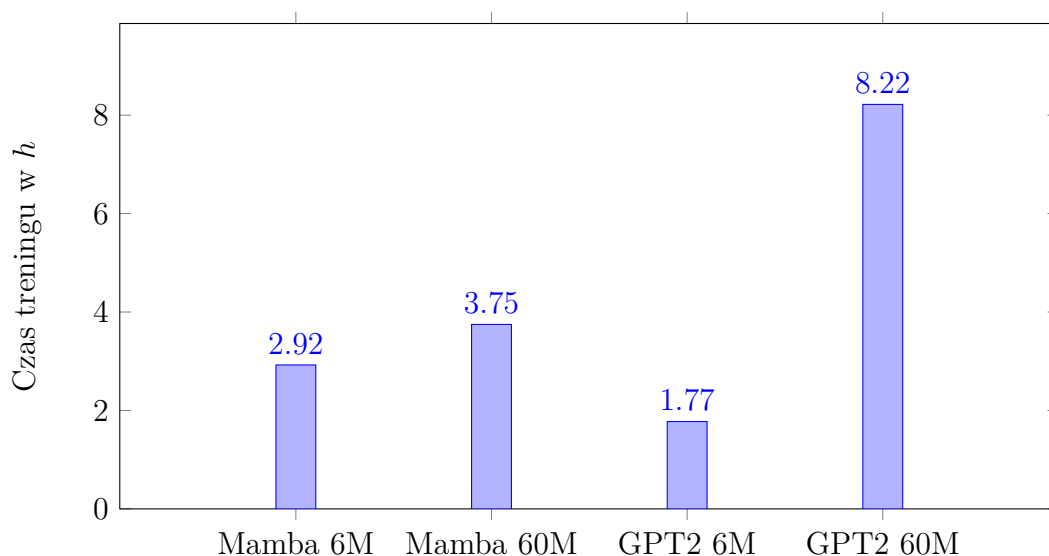
- `n_embed` – wymiarowość osadzeń oraz stanów ukrytych;
- `n_layer` – liczba warstw ukrytych w transformerze;
- `n_head` – liczba [złożeń algorytmu uwagi](#) w warstwach *multi-headed attention* [transformera](#);
- `n_inner` – wymiar środkowych warstw typu *feed-forward*.

Parametry modeli Mamba:

- `hidden_size` – wymiarowość osadzeń oraz stanów ukrytych;
- `state_size` – [rozmiar](#) stanów w modelu przestrzeni stanów;
- `num_hidden_layers` – ilość warstw ukrytych modelu.

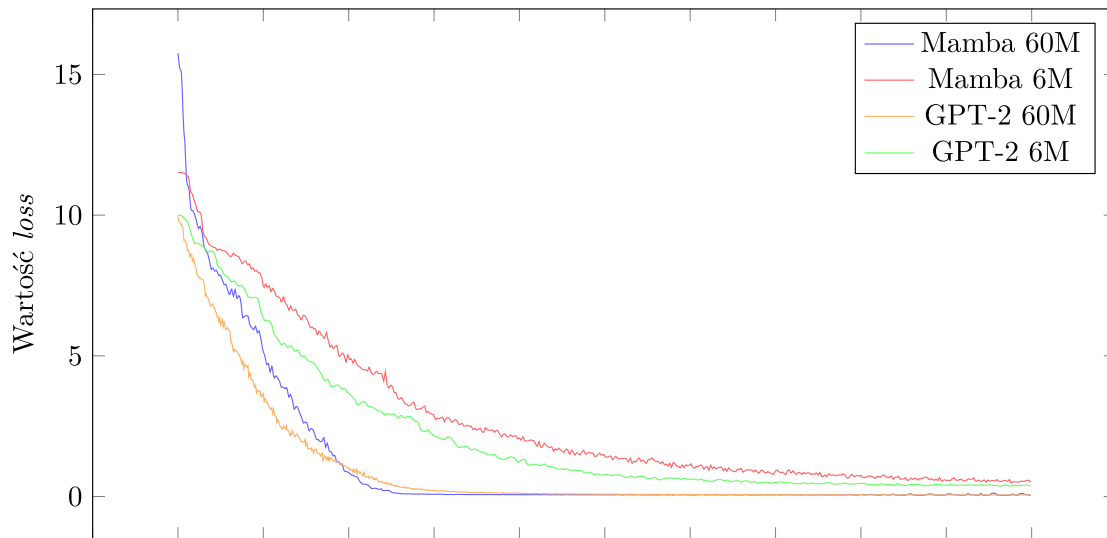
Pomimo tego, że istnieją pretrenowane modele generujące muzykę, w pracy został podjęty trening tego rodzaju modeli od samego początku. W związku z tym, że model Mamba jest dość nowym sposobem generowania sekwencji, nie powstały jeszcze modele muzyczne w tej architekturze, przez co porównywanie ich z pretrenowanymi modelami transformerowymi nie byłoby wiarygodnym testem.

Trening modeli odbywał się na karcie graficznej NVIDIA Tesla A100. Czas treningu modeli przedstawia wykres 3.1. Każdy model był trenowany przez około 800 tysięcy kroków, jednak ich dokładna liczba różniła się pomiędzy eksperymentami. Aby lepiej przedstawić różnice w czasie treningu, na wykresie przedstawiono tylko czas pierwszych 100 tysięcy kroków. Modele trenowane na plikach ABC wymagały zdecydowanie krótszego czasu treningu, co wynikało z mniejszej ilości unikatowych tokenów w stosunku do reprezentacji MIDI.



Rysunek 3.1.: Czas treningu 100 tysięcy kroków modeli.

Porównując czasy treningu obu architektur można dojść do następujących wniosków. Dość oczywistą obserwacją jest fakt, że trening modeli mniejszych jest krótszy niż tych większych. Pomimo tego, że model Mamba powinien trenować się szybciej niż modele transformerowe, nie jest to prawdą dla mniejszego modelu. Model transformerowy GPT2 6M trenował się o 39,4% szybciej niż jego odpowiednik w konkurencyjnej architekturze. Różnica w szybkości treningu jest zdecydowanie bardziej widoczna w dużych modelach, gdzie model transformerowy GPT2 60M trenował się prawie 120% dłużej niż model Mamba o tej samej liczbie parametrów. Na podstawie tych obserwacji można stwierdzić, że w kontekście dużych modeli architektura Mamba jest bardziej efektywna pod względem czasu treningu niż architektura GPT2.

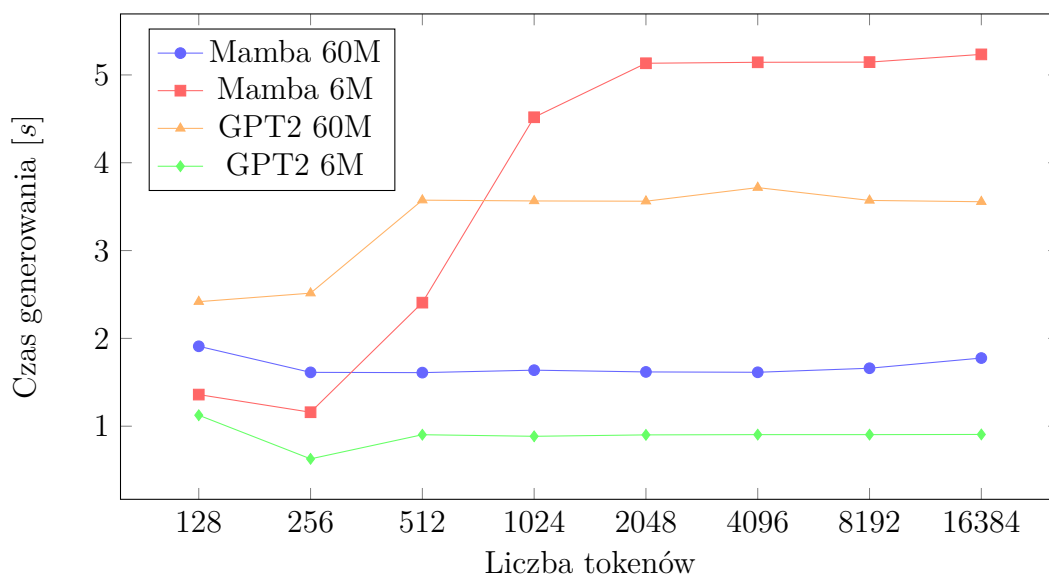


Rysunek 3.2.: Wykresy treningu pierwszych 250 tysięcy kroków modeli.

Na grafice 3.2 przedstawiono wartość funkcji straty podczas treningu modeli. W celu dokonania lepszej wizualizacji ograniczono ilość kroków do 250 tysięcy. Modele z większą liczbą parametrów (60M) szybciej osiągnęły niższą funkcję straty w porównaniu do modeli z mniejszą liczbą parametrów (6M). Dzieje się tak, ponieważ większe modele mają lepszą zdolność do reprezentowania skomplikowanych zależności w danych przy pomocy swoich parametrów, dzięki czemu mogą lepiej dopasować się do zbioru treningowego, co prowadzi do szybszej redukcji funkcji straty. Pomimo tego, że duże modele osiągnęły niższe wartości szybciej, wszystkie modele wytrenowały się do podobnego poziomu.

Na grafice 3.3 przedstawiono czas generowania konkretnej liczby tokenów przez modele. Najszybciej działa najmniejszy model GPT2, [który](#) dla każdej wartości generuje sekwencje poniżej jednej sekundy. Porównując duże modele, widać, że Mamba radzi sobie lepiej, szczególnie dla dłuższych sekwencji, działając prawie dwa razy szybciej. Dużym zaskoczeniem jest gwałtowny wzrost czasu dla modelu Mamba 6M.

Przyglądając się małemu oraz dużemu modelowi Mamby, można zauważyć, że różnica między nimi nie jest tak duża jak w przypadku modeli GPT2. Może wynikać to z tego, że model ten w swojej architekturze jest zoptymalizowany pod kątem obsługi bardzo dużej liczby parametrów. Najmniejszy model, jaki został wytrenowany i udostępniony przez twórców tej architektury, posiadał 128 milionów parametrów, co może sugerować, że dużo mniejsze modele nie zachowują obiecanych zalet związanych z ich osiąganiami.



Rysunek 3.3.: Czas generowania tokenów przez modele.

3.3. 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, [wygenerowane przykłady muzyczne zostaną przedstawione](#) w dwojaki sposób. W zależności od tego, czy model był wytrenowany na plikach MIDI [czy](#) 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 tym, że pliki MIDI są zapisem odegrania muzyki, a nie samej jej notacji, reprezentacja plików MIDI na pięciolinii potrafi nie być idealna, co jest spowodowane możliwymi nieregularnościami w rytmie.

Na grafice 3.4 przedstawiono wynik modelu GPT-2 60M w wersji ABC. Część oznaczoną kolorem niebieskim należy traktować jako [sekwencję](#) wejściową dla modelu, na której podstawie model generuje kolejne tokeny. Dodatkowo zostały podane kody kontrolne S, B oraz E. Reprezentację otrzymanej muzyki zapisanej na pięciolinii przedstawia grafika 3.5.

Należy zwrócić uwagę, że model został poproszony o wygenerowanie melodii w tonacji

C-dur (K:C), jednak, jak można zaobserwować na obu reprezentacjach, model zignorował tę część zapytania i wygenerował melodię w tonacji d-moll. 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 gamy), ponieważ takie fragmenty są najbardziej popularne i dominowały zbiór uczący.

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

Rysunek 3.4.: Wynik modelu.



Rysunek 3.5.: Wygenerowana melodia ABC modelem GPT-2 60M.

Przykład wielogłosowego utworu wygenerowanego przez model znajduje się na ilustracji 3.6. W tym przypadku model poprawnie utworzył model z zadanej tonacji D-dur. Model wygenerował utwór jedynie na podstawie kodów kontrolnych S:2, B:9, E:4, L:1/8, M:3/4, K:D podanych w zapytaniu dla modelu.



Rysunek 3.6.: Wielogłosowy fragment melodii wygenerowany modelem Mamba 60M.

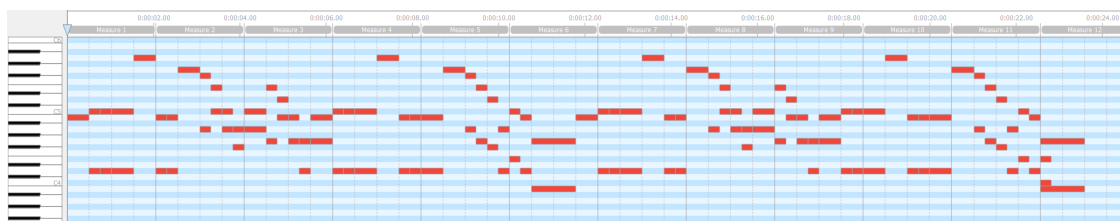
Generowanie plików MIDI jest procesem bardziej skomplikowanym niż generowanie muzyki w notacji ABC, w którym wystarczy przy pomocy klawiatury wprowadzić odpowiednie nuty. W tym przypadku predykcje są powiązane z tokenizatorem, który zamienia konkretne dźwięki na tokeny, które są rozumiane przez model. Najprostszym sposobem jest stworzenie początku sekwencji przy pomocy dowolnego edytora plików MIDI, a następnie otwarcie ich przy pomocy odpowiedniej biblioteki oraz zamianę na tokeny przy pomocy tego samego tokenizatora, który został użyty podczas treningu modelu.

Przykładowe predykcje stworzone przy pomocy modelu Mamba 6M można zobaczyć na grafikach 3.7 oraz 3.8. Fragment został wygenerowany jedynie na podstawie tokenu początku sekwencji, co sprawia, że model nie dostaje sugestii na temat tego, jaką melodię powinien stworzyć. Wygenerowany fragment zawiera dwa głosy grające w tym samym czasie.

Różnicą, która jest widoczna pomiędzy predykcjami 3.7 oraz 3.5, jest sposób generowania powtórzeń w muzyce. W notacji ABC odbywa się to przy pomocy znaków notacji muzycznej, które mówią o tym, jakie fragmenty mają zostać powtórzone. W przypadku predykcji plików MIDI, w drugiej linijce wygenerowanego fragmentu występuje powtórzenie pierwszej linijki (poza ostatnim taktem). Podobne powtórzenie mogłoby zostać zapisane na pięciolinii przy pomocy odpowiednich znaków. Tego rodzaju repetycja sekwencji wymaga od modelu zapamiętania jej w swoich wagach, a następnie przekopiowania jej w odpowiednie miejsce, co jest zadaniem nietrywialnym i podatnym na błędy.

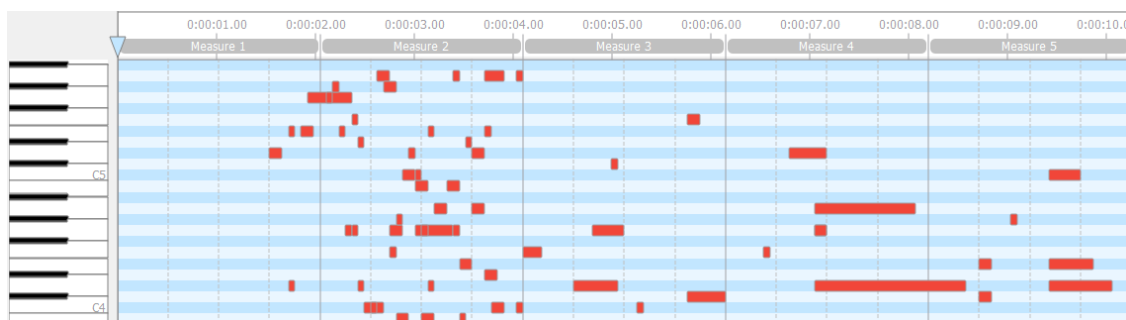


Rysunek 3.7.: Wygenerowana melodia MIDI zapisana w notacji muzycznej.



Rysunek 3.8.: Wygenerowana melodia MIDI modelem GPT2 6M.

Dotychczas w pracy zostały przedstawione fragmenty, które zostały wybrane przez autora jako wyjątkowo dobre i brzmiące wiarygodnie, jednak nie wszystkie predykcje należą do tej grupy. Część predykcji jest nieudanymi próbami generowania muzyki, która zwyczajnie nie brzmi dobrze. Tego rodzaju fragment przedstawiono na ilustracji 3.9. Zawiera on zdecydowanie za dużo nut grających w tym samym czasie, przez co prawdopodobnie nawet sprawny muzyk nie byłby go w stanie odegrać. Część fragmentów kończy się za szybko, ponieważ model dodał w nich znak końca sekwencji w niepoprawnym miejscu lub fragmenty zawierają zdecydowanie zbyt dużą liczbę pauz następujących po sobie, co wprowadza dość niezręczną ciszę w wygenerowanych fragmentach. Zdecydowanie większa część udanych generacji modelu powstała na podstawie plików ABC, co jest następstwem prostszej struktury tych plików. **Jednak również** w przypadku tego formatu, co jest szczególnie widoczne, kiedy model jest niedotrenowany, model jest w stanie wygenerować sekwencje, które nie są zgodne z notacją, przez co nie są one możliwe do odsłuchania w dostępnym oprogramowaniu. Takie fragmenty należy pomijać lub próbować poprawiać je ręcznie.



Rysunek 3.9.: Niepoprawny harmonicznie wygenerowany fragment MIDI.

3.4. Sposoby oceny wyników

Celem porównania otrzymanych wyników z dowolnych modeli jest ocena jakości ich pracy. W kontekście modeli dyskryminujących (ang. *discriminative models*), na przykład klasyfikatorów bądź modeli regresyjnych, ocena ich jakości jest dość prosta, ponieważ istnieje 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 generatywnych 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 **praktycznie niemożliwe**. Dla człowieka weryfikacja wyników modeli generatywnych, 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, na przykład w skali od 1-5, braną z zazwyczaj niewielkiej grupy ludzi. Niestety, metryka ta często nie jest wiarygodna ze względu na jej subiektywność oraz niewielką próbę badawczą. Aby wyeliminować te wady, serwisy, takie jak HuggingFace, udostępniają narzędzia dla członków społeczności. Narzędzia te 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*[39], 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 dobrać **akordy w celu harmonizacji** danej melodii. Niestety, większość takich narzędzi jest płatna i nie pozwala na zautomatyz-

zowaną analizę wielu plików. Dodatkowym problemem w analizie harmoniczej, szczególnie prowadzonej przez algorytmy, jest rozróżnienie harmonii wertykalnej oraz horyzontalnej. Rozróżnienie to zostało wytłumaczone przez Jacoba Colliera, kilkukrotnego laureata nagród *Grammy*, którego zdaniem akord, który nawet zagrany sam brzmi niepoprawnie, w odpowiednim kontekście i przez odpowiednią [progresję akordów w dalszej części muzyki](#), może mieć nadany sens, przez co cała sekwencja nabiera muzycznego piękna [40].

Brak obiektywnych metryk, którymi można się posłużyć podczas ewaluacji modelu, jest dodatkowym problemem w momencie próby dotrenowania 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”.

W zbiorze danych ABC kody kontrolne dodawane do zbioru danych mogą posłużyć jako przynajmniej częściowa weryfikacja tego, czy model wygenerował dany fragment na ich podstawie. Wyznaczenie ich dla wygenerowanego fragmentu nie odpowie na pytanie, czy dany fragment brzmi [harmonijnie](#). Podczas treningu model nie jest bezpośrednio penalizowany za to, że ich [nie uwzględnia](#). Jednak pomimo tego, większość wygenerowanych melodii trzyma się zapytania podanego przez użytkownika. W tabeli 3.3 porównano, ile z dziesięciu wygenerowanych sekwencji trzymało się kodów podanych przez użytkownika. W celu przeprowadzania eksperymentów jako zapytanie do modelu zostały przekazane wyłącznie kody kontrolne bez jakichkolwiek nut, co pozwoliło modelowi na maksymalną [swobodę](#) twórczą.

	liczba sekwencji	liczba taktów	podobieństwo fragmentów	metrum	tonacja
Mamba 6M	10	9	8	8	5
GPT-2 6M	10	10	7	9	4
Mamba 60M	10	10	9	10	7
GPT-2 60M	10	10	9	10	6

Tabela 3.3.: Porównanie wygenerowanych melodii ABC.

Duże modele lepiej poradziły sobie z odtworzeniem odpowiedniej tonacji, jednak różnice są niewystarczające, aby ocenić, czy używanie danej architektury ma [istotny wpływ](#) na jakość predykcji.

Kody kontrolne, dodane przed treningiem, mogły zostać policzone przy pomocy skryptów. Metrum oraz tonacja utworów niestety musiały zostać zweryfikowane ręcznie. Większość melodii zaczynała się przedtakterem, czyli taktem niepełnym, co wynikało z charakteru danych treningowych, które w większości składały się z tradycyjnych melodii, w których tego rodzaju element występuje dość często. Przedtakt to niepełny takt, który podkreśla to, że utwór rozpoczyna się od słabej, nieakcentowanej części przebiegu rytmicznego. Tego rodzaju zabieg muzyczny zazwyczaj powoduje, że ostatni takt jest skrócony o długość przedtaktu, co przy-

nosi trudności w algorytmicznym ustaleniu metrum utworu. Stworzenie utworu w konkretnej tonacji było najbardziej problematyczne dla modeli. Wynikało to z faktu, że w zbiorze danych nie każda tonacja była reprezentowana w takiej samej liczbie. Modele preferowały tonacje durowe, co również prawdopodobnie wynika z charakteru zbioru danych. Podczas eksperymentowania z parametrami modeli w czasie wykonywania pierwszych testów, dość jasno można było zauważyć, że modele bardzo dobrze radziły sobie z popularnymi tonacjami jak C-dur, G-dur lub d-moll. Mniej popularne tonacje, z większą ilością krzyżyków lub bemoli, takie jak es-moll lub Gis-dur, potrafiły być ignorowane, a melodia powstawała w prostszej tonacji. Takie błędy pojawiały się częściej w modelach o mniejszej liczbie parametrów, co można wytłumaczyć ich mniejszą złożonością, co przekłada się na gorsze predykcje.

W przypadku plików MIDI nie istnieje mechanizm podobny do kodów kontrolnych, na podstawie których można ocenić dokładność precykcji. Modele MIDI gorzej radziły sobie z generowaniem melodii, co wynika z bardziej złożonej struktury i dłuższych sekwencji wejściowych modelu, co było spowodowane przez potrzebę użycia tokenizatora. Melodie wygenerowane przez modele operujące na notacji ABC praktycznie zawsze były przyjemne do odsłuchania. W tabeli 3.4 zostały przedstawione subiektywne oceny autora dla dziesięciu wygenerowanych fragmentów MIDI. Podczas eksperymentu, każdy wygenerowany fragment powstał jedynie na podstawie tokenu rozpoczęcia sekwencji.

	poprawny	niepoprawny
Mamba 6M	4	6
GPT-2 6M	3	7
Mamba 60M	4	6
GPT-2 60M	6	4

Tabela 3.4.: Sybiektywna opinia autora na temat modeli MIDI.

Większe modele poradziły sobie lepiej od tych mniejszych, jednak aby otrzymać bardziej wiarygodne porównanie, należy przeprowadzić podobny eksperyment na większej populacji.

3.4.1. Testowanie melodii innymi LLM

Aby zweryfikować jakość predykcji, można zastosować podejście polegające na wytrenowaniu innego modelu na danych muzycznych. Tym sposobem powstanie model, który będzie w stanie ocenić oraz udzielić konstruktywnej krytyki na temat przedstawionego mu fragmentu muzycznego.

ChatMusician

ChatMusician to otwarty model językowy, który integruje zdolności muzyczne [41]. Model ten został zaprojektowany na bazie modelu LLaMA2 z siedmioma miliardami 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. Pytanie zadane do modelu brzmi następująco: „Wygenerowałem muzykę w notacji ABC. Czy mógłbyś ją przeanalizować i powiedzieć, czy brzmi naturalnie? W szczególności chciałbym prosić o opinię na temat ogólnej muzykalności, poczucia harmonii i tego, czy melodia i progresje akordów są spójne i przyjemne. Oto jej notacja:[...]”. W związku z tym, że model działa tylko dla języka angielskiego, należy przetłumaczyć pytanie na ten właśnie język.

Zapytanie o wygenerowany fragment zapisany na partyturze 3.5 przedstawiono [na rysunku 3.10](#).

```
PROMPT: I have some music in ABC notation that I generated. Could you
analyze it and tell me if it sounds natural? Specifically, I would like
your feedback on the overall musicality, the sense of harmony, and whether
the melody and chord progressions are coherent and pleasing. Here is the
ABC notation of the music:
X:1
L:1/8
M:3/4
K:C
D/2d/2d Ad fa | a/g/e/f/ g/f/e/d/ ^cA | d2 dA df | ag/f/ e^c A2 |
D2 d^c de | f2 ef/g/ a/g/f/e/ | d2 df e^c |1 d2 D4 :|2 d2 D2 fg |:
a2 g2 f2 | ed ^cd e/f/g | a2 g2 f2 | ed ^cd e/f g | a2 g2 f2 |
ed ^cd ef/g/ | a2 g2 fe |1 d2 de fg :|2 d2 D4 ||
```

Rysunek 3.10.: Pytanie dla modelu ChatMusician.

ANSWER: Overall Musicality: 1. The piece has a pleasant and flowing rhythm with a steady tempo. It features a mix of different instruments and voices, creating an interesting texture. The melodies are catchy and memorable, making the piece enjoyable to listen to. 2. Sense of Harmony: The harmonies in this piece are well-crafted and create a rich and full sound. The use of various chords and progressions adds depth and complexity to the music. The harmonic choices are appropriate for the style and genre of the piece. 3. Melody and Chord Progressions: The melody and chord progressions in this piece are coherent and pleasing. They follow a logical progression and fit together well. The melodic lines are well-written and expressive, adding to the emotional impact of the piece. The chord progressions provide a solid foundation for the melodies and add interest and variety to the composition. In conclusion, this piece of music has a strong sense of musicality, harmony, and coherence. The melodies and chord progressions are well-constructed and contribute to the overall appeal of the piece. It is a successful example of its genre and style.

Rysunek 3.11.: Odpowiedź modelu ChatMusician.

Model dość przychylnie spojrział na ten wygenerowany fragment, jednak jego odpowiedzi nie są specyficzne, o czym można się przekonać, czytając odpowiedź na grafice 3.11. Nie podejmuje on analizy konkretnych fragmentów. Dodatkowo udziela informacji o użytych wielu instrumentach, co nie występuje w przedstawionym mu fragmencie. W kontekście harmonii używa słów, takich jak *well-crafted*, *coherent*, *pleasing*. Słowa te są bardzo ogólne i nie stanowią dobrej analizy fragmentu muzycznego. Dodatkowo w opinii udzielonej przez model występują powtórzenia, które zmniejszają wiarygodność oceny. Model nie udzielił żadnej negatywnej odpowiedzi lub nie zasugerował możliwości zmiany.

Aby upewnić się, czy udzielanie wyłącznie pozytywnych informacji nie jest domyślnym zachowaniem modelu, należy poprosić model o ocenę fragmentu, który został specjalnie przygotowany, aby nie brzmiał dobrze dla ludzkiego ucha. Zapis tego fragmentu został zapisany na grafice 3.12. Dodatkowo, pomimo tego, że zapis jest poprawny w notacji ABC, niektóre takty zawierają więcej lub mniej dźwięków niż jest to przewidziane w metrum⁴. Odpowiedź modelu przedstawiono we fragmencie 3.13.

Obie opinie są bardzo do siebie podobne. Okazuje się, że model nie jest w stanie dobrze rozpoznać to, czy muzyka jest rzeczywiście przyjemna dla ucha czy nie. Nie wychwycił problemów z rytmem, który nie jest jednakowy dla każdego taktu oraz nie wyraził żadnej konstruktywnej krytyki. Niestety, model *ChatMusician* nie jest dobrym [rozwiązaniem problemu krytycznej oceny wygenerowanej muzyki](#).



Rysunek 3.12.: Specjalnie stworzony niepoprawny fragment muzyki.

ANSWER: Overall Musicality: The piece has a pleasant and flowing quality to it, with a mix of different rhythms and tempos. The use of dynamics and articulations adds depth and interest to the composition. 2. Sense of Harmony: The harmonies used in this piece are mostly consonant and create a sense of stability and balance. There are occasional dissonances that add tension and color to the music. 3. Melody and Chord Progressions: The melodies in this piece are well-crafted and follow a logical progression. They are easy to sing or play and have a clear sense of direction. The chord progressions are also coherent and support the melodies effectively. In conclusion, this piece has a good overall musicality, sense of harmony, and coherence between its melody and chord progressions. It is an enjoyable and engaging composition.

Rysunek 3.13.: Odpowiedź modelu *ChatMusician* na fałszywy fragment.

GPT-4

Aby porównać możliwości oceny modeli LLM, należy użyć silniejszego modelu. GPT-4 to model, który obecnie najlepiej wypada na praktycznie każdym benchmarku [42]. W porównaniu do swojego poprzednika, GPT-3, charakteryzuje się większą liczbą parametrów, co pozwala mu na lepsze rozumienie i generowanie skomplikowanych tekstów. GPT-4 potrafi również wykazywać większą spójność i koherencję w dłuższych wypowiedziach, lepiej interpretować kontekst oraz odpowiadać na bardziej złożone pytania. W związku z tym, że model nie jest *open source*, jego *dokładna* konfiguracja nie jest znana, jednak przypuszcza się, że ma 1,67 trylionów parametrów [43]. Nie był on trenowany bezpośrednio na danych muzycznych, jednak poprzez swój rozmiar, który jest potencjalnie ponad 2000 razy większy od modelu *ChatMusician*, powinien posiadać przynajmniej podstawową znajomość *zagadnień muzycznych*. Wadą tego rozwiązania jest fakt, że GPT-4 jest dostępny jedynie po wykupieniu

3. Przeprowadzenie eksperymentów

subskrypcji na stronie OpenAI, która kosztuje 20\$ na miesiąc.

Do modelu zostało przekazane takie samo pytanie jak w wyżej opisanym przypadku.

```
To evaluate the musicality and harmony of your ABC notation, I will provide an analysis based on the key, melody, rhythm, and chord progressions. Here is a detailed breakdown of your piece:

### Key and Meter
- **Key:** C major
- **Meter:** 3/4 (Waltz time)

### Melody and Harmony
1. **Opening Phrase:**
   - `D/2d/2d Ad fa | a/g/e/f/ g/f/e/d/ ^cA`
   - This phrase introduces a mix of stepwise motion and leaps, giving a dynamic start. The alternation between the upper and lower notes creates a call-and-response effect. The use of the raised 4th degree (^c) adds a touch of chromaticism.

2. **Middle Phrase:**
   - `d2 dA df | ag/f/ e^c A2`
   - The repetition of 'd' with variations in intervals suggests a thematic consistency. The leap to 'f' and back to 'd' maintains interest. The movement towards 'A' anchors the phrase back to the tonic's dominant, providing harmonic stability.

3. **Transition:**
   - `D2 d^c de | f2 ef/g/ a/g/f/e/`
   - The use of the D minor chord (D-F-A) transitioning to E minor (E-G-B) and then back to fragments of D minor adds tension and resolution.

4. **Closing Phrase:**
   - `d2 df e^c |1 d2 D4 :|2 d2 D2 fg`
   - This phrase returns to the tonic D and introduces a cadential figure (D-F-E-^C). The first ending resolves to the tonic D major chord, while the second ending moves into a continuation with 'fg', creating a sense of forward motion.

### Second Part
1. **Opening Phrase:**
   - `a2 g2 f2 | ed ^cd e/f/g`
   - This section starts with a descending sequence, providing contrast to the first part. The melodic minor descent with a chromatic touch (e-^c-d) maintains interest.

2. **Middle Phrase:**
   - `a2 g2 f2 | ed ^cd e/f g`
   - The repetition reinforces the theme while subtle variations in articulation (slurs) add expressiveness.

3. **Transition and Closing:**
   - `a2 g2 f2 | ed ^cd ef/g/ | a2 g2 fe`
   - The alternation between stepwise motion and leaps ensures coherence and variety. The repeated 'a-g-f' figure creates familiarity, while 'e-^c-d ef/g/' keeps it fresh.

4. **Final Cadence:**
   - `1 d2 de fg :|2 d2 D4 ||`
   - The final cadence resolves nicely with the repeating 'de fg' figure leading to the tonic 'D' again, ensuring a satisfying conclusion.

### Overall Musicality
- **Melodic Coherence:** The melody has a clear direction with a good mix of stepwise motion and leaps, maintaining interest.
- **Harmony:** The harmony is mostly implied through the melodic line with occasional chromatic touches adding tension and resolution. The frequent returns to tonic and dominant notes (C and G) provide a sense of harmonic stability.
- **Rhythmic Flow:** The piece maintains a consistent waltz rhythm (3/4 time), which gives it a dance-like quality, adding to its musicality.

### Conclusion
The piece sounds natural and coherent, with a pleasing melody and harmonic progression. The use of chromatic notes and variations in articulation adds interest and depth. It effectively balances repetition with variation, ensuring both familiarity and novelty for the listener.
```

Rysunek 3.14.: Odpowiedź modelu GPT-4 na wygenerowany fragment.

Jak można zauważyć na grafice 3.14, ocena fragmentu muzyki przez GPT-4 jest dość długa oraz bardzo dokładna. Podczas analizy GPT-4 rozbija muzykę na główne fragmenty, a następnie każdy z tych fragmentów poddaje szczegółowej analizie, rozdzielając je dalej na początek, środek, przejście i zakończenie. Dla każdej z wyróżnionych części udziela [odpowiedzi opartej o elementy analizy harmonicznej](#), mówiącej o progresji akordów, używając technicznego słownictwa jak dominanta lub tonika. GPT-4 rozpoznaje kierunkowość muzyki,

czyli to, czy melodia jest opadająca czy rosnąca. Model jest w stanie wydać opinię na temat rytmu oraz tego, w jaki sposób melodia komponuje się z rytmem. Przez to, że tonacja melodii została oznaczona jako C-dur, a tak naprawdę utwór jest w d-moll, model nie do końca jest pewny, przez co udziela niepoprawnej informacji o dźwięku cis jako czwartym stopniu gamy. Pomimo tego, pod koniec wypowiedzi, została udzielona prawidłowa informacja o zakończeniu utworu na tonice, czyli dźwięku D. Na końcu wypowiedzi model udziela podsumowania. Dodatkową zaletą jest fakt, że odpowiedź jest udzielona w formacie *markdown*, przez co tekst jest podzielony na logiczne fragmenty, a kluczowe słowa są oznaczone pogrubioną czcionką.

Aby eksperyment był wiarygodny, należy również przedstawić modelowi fragment „złej” muzyki. Odpowiedź modelu przedstawiono na ilustracji 3.15.

```
Based on the ABC notation provided, here is an analysis of the music in terms of musicality, harmony, and melody coherence:

### Overall Musicality
The notation seems to follow a structure that could be pleasing. The use of different note values and ornamentations like the grace note `{g}` add interest to the melody. The notation implies a mix of staccato and legato sections, which can contribute to a dynamic musical flow.

### Harmony
The harmony is somewhat implied rather than explicitly stated. The melody suggests movement between different tonal centers, but the use of accidentals (like `~g`, `_A`, `^c`, etc.) makes the harmony more complex. It moves away from a strictly diatonic C major scale, introducing chromaticism which can add color but might sound unusual without proper context or accompanying chords.

### Melody Coherence
The melody uses a variety of intervals and note values, which gives it a sense of motion. However, the frequent use of accidentals (`~g`, `_A`, `^c`, `_A`, `^B`) can make it sound less coherent if not handled carefully. These accidentals can create tension and release, which is good for musical interest but must be resolved correctly to sound pleasing.

### Specific Sections
- **Measure 1**:: Starts with a simple C major outline but immediately introduces accidentals (`~g`, `_A`), which suggest movement towards different tonal areas.
- **Measure 2**:: The grace note `{g}` and the intervals used (`D1 f2 ^c2`) create a surprising twist.
- **Measure 3**:: The use of `G2 e1/2 ^g2 _A2` suggests a continuation of the chromatic movement, making it intriguing.
- **Measure 4**:: The abrupt stop with `B3` could serve as a point of tension or punctuation.
- **Measure 5**:: A return to `C2 e2 ^g2 [_A2C]` reintroduces some stability but with the added complexity of accidentals.
- **Measure 6**:: `B3 D22` maintains the movement and suggests a possible resolution.
- **Measure 7**:: The phrase `C2 _e2 ^g2 _A2` is quite chromatic and could be disorienting without a clear harmonic context.
- **Measure 8**:: The closing `B3 D1 f2 ^c2` should ideally resolve the chromatic tension but continues to add accidentals, leaving a sense of unresolved tension.

### Conclusion
The melody and chord progressions can be coherent and pleasing if the chromaticism and accidentals are used intentionally to create tension and resolution. It may sound natural in a modern or experimental context but could be challenging in a more traditional tonal framework. To improve musicality, ensure that the chromatic notes resolve in a way that feels intentional, and consider adding harmonic support to clarify the underlying chord progressions.
```

Rysunek 3.15.: Odpowiedź modelu GPT-4 na fałszywy fragment.

W przypadku niepoprawnego fragmentu muzycznego model w dość niepewny sposób udziela odpowiedzi. Często używa takich wyrażen, jak „struktura **może** być przyjemna”, „...**może** wpływać na dynamikę”, „**niespodziewany zwrot** muzyki”. Sugeruje to, że model rozumie dość nienaturalną strukturę przekazanej mu melodii. Dokonuje również analizy konkretnych taktów, gdzie dość krytycznie wypowiada się na ich temat. Oferuje również konstruktywną krytykę, zauważając, że częste stosowanie przypadkowych dźwięków i fraz chromatycznych może wywoływać poczucie dezorientacji i nierozwiązanego napięcia. Gwałtowne zatrzymania i złożone interwały są doceniane za to, że mogą sprawić, że dany fragment będzie bardziej interesujący, ale są również podkreślane jako punkty, które mogłyby skorzy-

stać z wyraźniejszego harmonicznego rozwiązania. Ogólnie rzecz ujmując, modelowi udało się zwrócić uwagę na szczególne fragmenty w melodii i poddać je krytycznej analizie, która nie była wyłącznie przychylna.

Okazało się, że model GPT-4, pomimo tego, że nie był specjalnie wyszkolony do tego celu, lepiej przeanalizował muzykę, wykazując zdolność do precyzyjnej identyfikacji i oceny harmonicznych struktur w utworach. Odpowiedź modelu na temat specjalnie przygotowanego błędnego fragmentu może zostać portaktowana jako krytyczna, jednak zawiera dość dużą ilość informacji pozytywnych. Prawdopodobnie wynika to z charakteru olbrzymich zbiorów danych treningowych, w skład których mogą wchodzić recenzje muzyczne, które zwykle nie są pisane na temat niepoprawnej muzyki. Tego rodzaju uprzedzenie może prowadzić do częstszego generowania pozytywnych recenzji. Model *ChatMusician*, który również przeanalizował te same fragmenty, przedstawił bardzo podobną ocenę dla obu, zarówno poprawnego, jak i niepoprawnego harmonicznie, co wskazywało na jego ogólną analizę, nieodróżniającą istotnych niuansów harmonicznych.

Istnieje technika nazwana RLHF (*reinforcement learning with human feedback*), która jest wykorzystywana w trenowaniu lub dotrenowaniu modeli sztucznej inteligencji [44]. Polega ona na zbieraniu danych od ludzi na temat działania modelu i wykorzystaniu tych danych do poprawy jego wydajności. W praktyce wygląda to tak, że model generuje odpowiedzi na różne zapytania, a ludzie oceniają te odpowiedzi pod kątem ich poprawności, użyteczności, spójności czy innych kryteriów jakościowych. Na podstawie tych ocen model jest wzmacniany, aby generować coraz lepsze odpowiedzi w przyszłości. Przedstawione wyniki analizy modelu GPT-4 mogą być brane pod uwagę jako dobra opinia na temat tego, czy dany wygenerowany fragment jest poprawny czy nie. Tym sposobem istnieje szansa na rozwój tego rodzaju sposobu dotrenowywania modeli generatywnych. Problemem jest obecnie koszt używania modelu od OpenAI, który posiada codzienne limity, które, co prawda, można zwiększyć za określoną kwotę. Niestety, nie jest to obecnie dobrze skalujące się rozwiązanie, szczególnie dla użytku niekomercyjnego. Mimo wszystko jest to możliwy sposób dotrenowywania modeli generatywnych muzyki, pod warunkiem odpowiedniej inżynierii zapytania przekazywanego do modelu nauczyciela.

3.4.2. Muzyczny test Turinga

Test Turinga dla muzyki, analogicznie do klasycznego testu Turinga, jest próbą oceny czy maszyna może wykazać się poziomem inteligencji muzycznej, który jest nieodróżnialny od ludzkiego. Pomysł ten pojawił się po raz pierwszy w 1988 roku w *Computer Music Journal* jako forma identyfikacji inteligencji maszynowej w muzyce [45]. Test ten polega na tym, że maszyna i człowiek uczestniczą w sesji improwizacyjnej (*jam session*), a zadaniem oceniającego jest rozpoznanie, który z uczestników jest maszyną. Maszyna musi być w stanie wykryć rytm, rozpoznać, co gra człowiek i odpowiedzieć na to w czasie rzeczywistym, tworząc improwizowane „solówki”, które swoją strukturą będą niejako „odpowiedzią” na to, co zagrał

człowiek.

Oryginalny test Turinga, zaproponowany przez Alana Turinga w 1950 roku, polegał na tym, że sędzia prowadził rozmowę z dwoma podmiotami, jednym ludzkim i jednym maszynowym, bez wiedzy, który jest który [46]. Jeśli sędzia nie mógłby odróżnić maszyny od człowieka na podstawie odpowiedzi, maszyna przechodziła test, wykazując inteligencję na poziomie ludzkim. Test Turinga dla muzyki rozszerza tę koncepcję na kontekst muzyczny, gdzie zamiast rozmowy mamy interaktywną improwizację muzyczną. Kluczowym elementem jest tutaj nie tylko wynik, ale również proces tworzenia muzyki w czasie rzeczywistym, który musi być na tyle ludzki, aby oceniający nie mógł odróżnić maszyny od człowieka. Obecnie każdy model, który generuje muzykę, skupia się na wyniku, a nie na procesie. Sztuczna inteligencja musiałaby być w stanie reagować na zmiany w muzyce w czasie rzeczywistym, co wymaga nie tylko analizy dźwięku, ale również interpretacji wizualnych sygnałów i dopasowania do tempa i stylu gry człowieka. Jest to znacznie bardziej skomplikowane niż generowanie statycznych utworów muzycznych. Czynniki te sprawiają, że chociaż modele generatywne są w stanie generować muzykę, przejście muzycznego testu Turinga, który wymaga głębokiej interakcji i ludzkiego zrozumienia procesu twórczego, pozostaje wyzwaniem trudnym do osiągnięcia. Niemniej jednak możliwe jest poddanie wygenerowanych fragmentów bezpośredniej ocenie ludzkiej. W tym celu została przygotowana playlista w serwisie YouTube zawierająca kilka wygenerowanych fragmentów. W ten sposób każdy może przekonać się sam, czy wygenerowane melodie mogą konkurować z tymi, stworzonymi przez człowieka [47].

4. Zakończenie

Omawiane w pracy modele przeznaczone do generowania i analizy tekstu wykazały obiecujące zdolności również w analogicznych zastosowaniach w kontekście muzyki. Uznane modele transformerowe dowiodły możliwość radzenia sobie z tymi zagadnieniami. Poprzez intensywne badania w tematach NLP pojawiają się nowe rozwiązania, na przykład modele Mamba, które dokonują generowania prostej muzyki w jakości zbliżonej do tej pochodzącej z transformerów. Poruszone porównanie zapisu ABC oraz MIDI pozwala wybrać pasujące dane i model do potrzeb użytkownika, aby otrzymać rozwiązanie, które będzie dla niego najbardziej zadowalające.

Przedstawione metody generowania muzyki mogą być dalej rozwijane głównie poprzez zwiększanie liczby parametrów modeli. Im większy model, tym jest on w stanie generować bardziej skomplikowane oraz innowacyjne melodie. Należy zwrócić uwagę również na zróżnicowanie danych treningowych. Ze względu na dużą liczbę melodii irlandzkich w zbiorze danych, wygenerowane kompozycje często posiadają charakterystyczne irlandzkie brzmienie. Większe modele mogą nadmiernie dopasować się do przedstawionych im danych, przez co dywersyfikacja oraz gromadzenie większej ilości melodii jest istotnym elementem w kontekście poprawienia jakości modelu. Istniejące zbiory danych ABC można rozbudować o dodatkowe kody kontrolne, które, na przykład mogą kodować ilość zapisanych głosów we fragmentach muzycznych. Tego rodzaju dodatek zwiększy kontrolę użytkownika nad procesem generowania muzyki. Dodatkowo praca poruszyła dość innowacyjne rozwiązanie oceny muzyki przy pomocy LLM. Rozwiązanie to może zostać potraktowane jako załączek sposobu końcowego dotrenowywania parametrów modelu przy pomocy technik uczenia ze wzmacnianiem, gdzie to nie człowiek, a algorytm dostarcza informację zwrotną na temat wygenerowanych fragmentów.

Przedstawione rozwiązania mogą prowadzić do prób stworzenia plug-inów lub podobnych narzędzi do programów, takich jak *FL Studio* lub *MuseScore*, które pozwalałyby na generatywne kończenie muzyki, która jest tworzona przez kompozytora. Podobne rozwiązania zostały wprowadzone przez Adobe w programie Photoshop, gdzie osoba mogła zaznaczyć dany fragment grafiki, a następnie poprosić AI, aby dokonało pewnych korekt według polecenia. Można wyobrazić sobie, że na przykład we wcześniej wymienionym *FL Studio*, artysta tworzy fragment muzyki, a następnie dostaje kilka odpowiedzi, w jaki sposób mógłby wzbogacić lub kontynuować utwór. W taki sam sposób osoba tworząca zapis muzyczny w *MuseScore*, mogłaby dostawać propozycje od modelu, który dostaje automatycznie wygenerowaną notację muzyczną w formacie ABC na podstawie innego zapisu. Tym sposobem, twórcy mo-

gliby wspierać swoją kreatywność przy pomocy metod sztucznej inteligencji. Przedstawione w pracy badania należy traktować jako początek możliwości dalszego rozwoju modeli opartych o architekturę LLM w kontekście muzyki. Otrzymane wyniki zachęcają do dalszych eksperymentów, na podstawie których mogą powstać bardziej zaawansowane modele generatywne dla muzyki.

5. Podziękowania

Dziękujemy Polskiej Infrastrukturze Komputerów wysokiej wydajności PLGrid (Centrum HPC: ACK Cyfronet AGH) za udostępnienie sprzętu komputerowego i wsparcie w ramach grantu obliczeniowego nr PLG/2024/017281 oraz PLG/2024/017191.

We gratefully acknowledge Polish high-performance computing infrastructure PLGrid (HPC Center: ACK Cyfronet AGH) for providing computer facilities and support within computational grant no. PLG/2024/017281 and PLG/2024/017191.

Bibliografia

- [1] Daniel Adiwardana i in. *Towards a Human-like Open-Domain Chatbot*. 2020. arXiv: 2001.09977 [cs.CL].
- [2] Anton Chernyavskiy, Dmitry Ilvovsky i Preslav Nakov. „Transformers: “The End of History” for Natural Language Processing?” W: *Machine Learning and Knowledge Discovery in Databases. Research Track*. Red. Nuria Oliver i in. Springer International Publishing, 2021, s. 677–693. ISBN: 978-3-030-86523-8. DOI: 10.1007/978-3-030-86523-8.
- [3] Örjan Sandred, Mikael Laurson i Mika Kuuskankare. „Revisiting the Illiac Suite - A rule-based approach to stochastic processes”. W: *Sonic Ideas/Ideas Sonicas 2* (sty. 2009), s. 42–46.
- [4] Google AI. *Magenta*. URL: <https://magenta.tensorflow.org/> (term. wiz. 25.06.2024).
- [5] Adam Roberts i in. „A Hierarchical Latent Vector Model for Learning Long-Term Structure in Music”. W: *Proceedings of the 35th International Conference on Machine Learning*. Red. Jennifer Dy i Andreas Krause. T. 80. Proceedings of Machine Learning Research. PMLR, 2018, s. 4364–4373.
- [6] Jesse Engel i in. „GANSynth: Adversarial Neural Audio Synthesis”. W: *International Conference on Learning Representations*. 2019.
- [7] Ashish Vaswani i in. „Attention is All you Need”. W: *Advances in Neural Information Processing Systems*. Red. I. Guyon i in. T. 30. Curran Associates, Inc., 2017.
- [8] Cheng-Zhi Anna Huang i in. *Music Transformer*. 2018. arXiv: 1809.04281 [cs.LG].
- [9] Yueyue Zhu i in. *A Survey of AI Music Generation Tools and Models*. 2023. arXiv: 2308.12982 [cs.SD].
- [10] Botao Yu i in. „Museformer: Transformer with Fine- and Coarse-Grained Attention for Music Generation”. W: *Advances in Neural Information Processing Systems*. Red. S. Koyejo i in. T. 35. Curran Associates, Inc., 2022, s. 1376–1388.
- [11] Andrea Agostinelli i in. *MusicLM: Generating Music From Text*. 2023. arXiv: 2301.11325 [cs.SD].
- [12] *Multimedia Programming Interface and Data Specifications 1.0*. URL: <https://www.aelius.com/njh/wavemetatools/doc/riffmci.pdf> (term. wiz. 25.06.2024).

- [13] The International MIDI Association. *Standard MIDI-File Format Spec. 1.1*. URL: <http://www.music.mcgill.ca/~ich/classes/mumt306/midiformat.pdf> (term. wiz. 25.06.2024).
- [14] Chris Walshaw. *A brief history of ABC*. URL: <https://abcnotation.com/history> (term. wiz. 25.06.2024).
- [15] Steve Allen. *Beethoven Symphony No. 7, Movement 2 in ABC*. URL: <https://www.ucolick.org/~sla/abcmusic/sym7mov2.html> (term. wiz. 25.06.2024).
- [16] Jonathan J. Webster i Chunyu Kit. „Tokenization as the Initial Phase in NLP”. W: *Proceedings of the 14th Conference on Computational Linguistics - Volume 4*. COLING '92. Nantes, France: Association for Computational Linguistics, 1992, s. 1106–1110. DOI: 10.3115/992424.992434.
- [17] 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.
- [18] Shangda Wu i in. *TunesFormer: Forming Irish Tunes with Control Codes by Bar Patching*. 2023. arXiv: 2301.02884 [cs.SD].
- [19] Nitish Shirish Keskar i in. „CTRL - A Conditional Transformer Language Model for Controllable Generation”. W: *arXiv preprint arXiv:1909.05858* (2019).
- [20] Darrell Conklin. *Bach Chorales*. UCI Machine Learning Repository. DOI: 10.24432/C5GC7P.
- [21] I. S. Duff, A. M. Erisman i J. K. Reid. *Direct Methods for Sparse Matrices*. Oxford University Press, sty. 2017. ISBN: 9780198508380. DOI: 10.1093/acprof:oso/9780198508380.001.0001.
- [22] Curtis Hawthorne i in. „Enabling Factorized Piano Music Modeling and Generation with the MAESTRO Dataset”. W: *International Conference on Learning Representations*. 2019.
- [23] 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.
- [24] Ian Goodfellow, Yoshua Bengio i Aaron Courville. *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press, 2016, s. 409–411.

-
- [25] Sergey Ioffe i Christian Szegedy. „Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift”. W: *Proceedings of the 32nd International Conference on Machine Learning*. Red. Francis Bach i David Blei. T. 37. Proceedings of Machine Learning Research. Lille, France: PMLR, lip. 2015, s. 448–456.
- [26] Sepp Hochreiter i Jürgen Schmidhuber. „Long Short-term Memory”. W: *Neural computation* 9 (grud. 1997), s. 1735–80. DOI: 10.1162/neco.1997.9.8.1735.
- [27] Christopher Olah. *Understanding LSTM Networks*. 27 sierp. 2015. URL: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/> (term. wiz. 25.06.2024).
- [28] Sepp Hochreiter. „The Vanishing Gradient Problem During Learning Recurrent Neural Nets and Problem Solutions”. W: *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* 6 (kw. 1998), s. 107–116. DOI: 10.1142/S0218488598000094.
- [29] Zhiyong Cui i in. *Deep Bidirectional and Unidirectional LSTM Recurrent Neural Network for Network-wide Traffic Speed Prediction*. 2019. arXiv: 1801.02143 [cs.LG].
- [30] Jacob Devlin i in. „BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”. W: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Red. Jill Burstein, Christy Doran i Tamar Solorio. Minneapolis, Minnesota: Association for Computational Linguistics, czer. 2019, s. 4171–4186. DOI: 10.18653/v1/N19-1423.
- [31] Maarten Grootendorst. *A Visual Guide to Mamba and State Space Models*. 19 lut. 2024. URL: <https://newsletter.maartengrootendorst.com/p/a-visual-guide-to-mamba-and-state> (term. wiz. 25.06.2024).
- [32] 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.
- [33] 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.
- [34] Albert Gu i in. „HiPPO: Recurrent Memory with Optimal Polynomial Projections”. W: *Advances in Neural Information Processing Systems*. Red. H. Larochelle i in. T. 33. Curran Associates, Inc., 2020, s. 1474–1487.
- [35] Albert Gu i Tri Dao. *Mamba: Linear-Time Sequence Modeling with Selective State Spaces*. 2023. arXiv: 2312.00752 [cs.LG].
- [36] Philip Gage. „A new algorithm for data compression”. W: *The C Users Journal* 12.2 (lut. 1994), s. 23–38. ISSN: 0898-9788.
-

- [37] Nathan Fradet i in. *Byte Pair Encoding for Symbolic Music*. 2023. arXiv: 2301.11975 [cs.LG].
- [38] Alec Radford i in. *Language Models are Unsupervised Multitask Learners*. 2019. URL: <https://openai.com/index/better-language-models/> (term. wiz. 24.06.2024).
- [39] *The TTS Arena*. URL: <https://huggingface.co/blog/arena-tts> (term. wiz. 25.06.2024).
- [40] Jacob Collier. *That's not a wrong note, you just lack confidence*. URL: https://www.youtube.com/watch?v=meha_FCcHbo (term. wiz. 25.06.2024).
- [41] Ruibin Yuan i in. *ChatMusician: Understanding and Generating Music Intrinsically with LLM*. 2024. arXiv: 2402.16153 [cs.SD].
- [42] OpenAI i in. *GPT-4 Technical Report*. 2024. arXiv: 2303.08774 [cs.CL].
- [43] Maximilian Schreiner. *GPT-4 architecture, datasets, costs and more leaked*. 11 lip. 2023. URL: <https://the-decoder.com/gpt-4-architecture-datasets-costs-and-more-leaked/> (term. wiz. 25.06.2024).
- [44] Long Ouyang i in. „Training language models to follow instructions with human feedback”. W: *Advances in Neural Information Processing Systems*. Red. S. Koyejo i in. T. 35. Curran Associates, Inc., 2022, s. 27730–27744.
- [45] 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>.
- [46] A. M. TURING. „I.—COMPUTING MACHINERY AND INTELLIGENCE”. W: *Mind* LIX.236 (paź. 1950), s. 433–460. ISSN: 0026-4423. DOI: 10.1093/mind/LIX.236.433.
- [47] Filip Ręka. URL: https://www.youtube.com/playlist?list=PLT1HvMf6Qr_4wWwesELL9FQ9A1edJqT6c (term. wiz. 25.06.2024).