2025年8月 MUS研究会 音楽音響処理(1)

音楽基盤モデルは 音高情報を螺旋構造に埋め込むか?

八木 颯斗 1 ,高道 慎之介 1,2

(1: 慶大, 2: 東大)



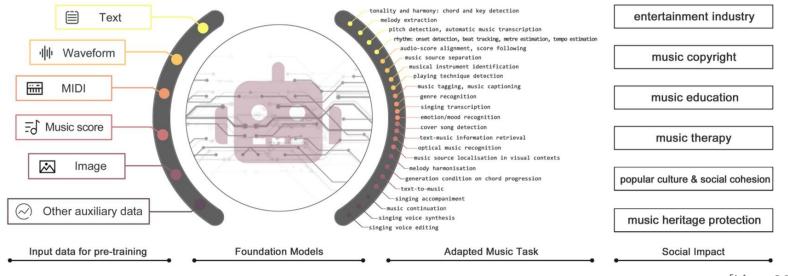
研究背景

基盤モデルとは

音楽基盤モデルは音高情報を螺旋構造に埋め込むか?

特定のタスクだけに特化して訓練されたモデルではなく 大規模な自己教師あり学習によって事前学習されることで多数の下流タスクに 対応可能な汎用的な機械学習モデル

例)BERT,GPT,Stable Diffusion,Jukebox,MusicGenなど



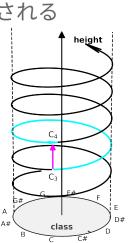
音高の表現

音楽基盤モデルは**音高情報**を**螺旋構造**に埋め込むか?

- 音高の知覚的表現は以下の2つの心理的属性から構成される
 - ピッチクラス (pitch class)
 - トーンクロマとも呼ばれ、音高の循環的な次元を指す
 - ▸ 基本的に12のピッチクラスとして定義
 - ピッチハイト (pitch height)
 - 音高の垂直的な次元であり、音を「低い」から「高い」へと順序付ける
 - 楽譜上では、どのオクターブに属するかによってこの高さが示される
 - 例) C4とC5は同じ「ド」の音でも,1オクターブの違いを表す

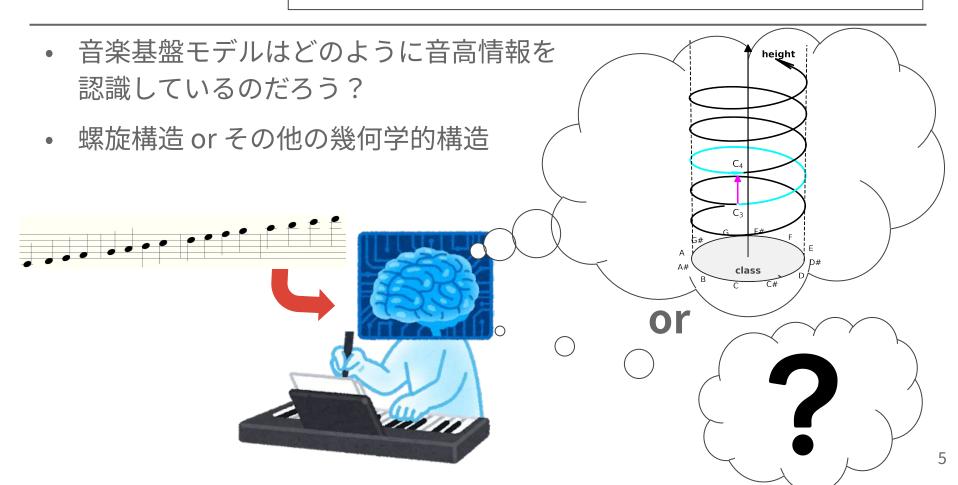


• 音高は螺旋モデルで説明可能 [Sheard, 1982]

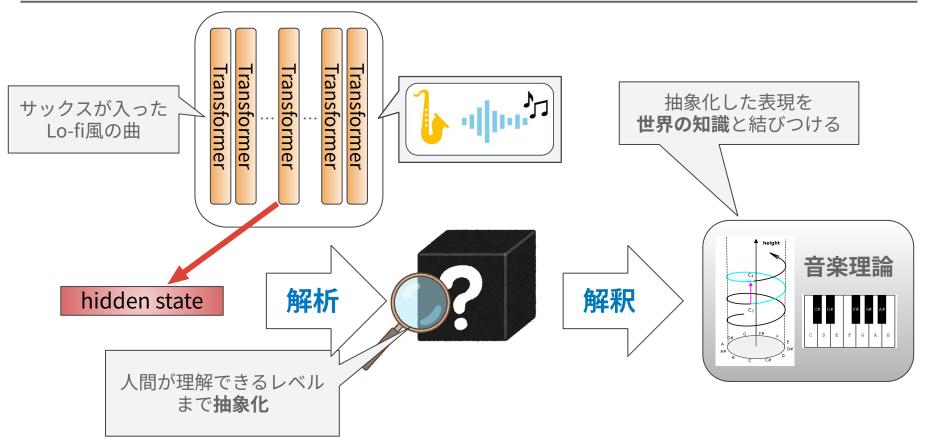


音高の表現

音楽基盤モデルは音高情報を螺旋構造に埋め込むか?



内部解析と解釈



何が嬉しい?

- 透明性の確保
 - モデルが「本当に理解しているか」検証できる
 - 音高に着目 → 人間の知覚に近い構造を獲得しているかどうかを検証
- 音楽の知見を発展
 - 音楽の仮説を検証する手段になる可能性
 - 人間のピッチ知覚の真の構造を知るヒントになる
- 制御性の向上
 - 音楽的要素を特定し、細かく直感的に制御可能
 - 創作活動において欠かせないユーザーの主体性を確保に繋がる

関連研究

本実験で使用する音楽基盤モデル

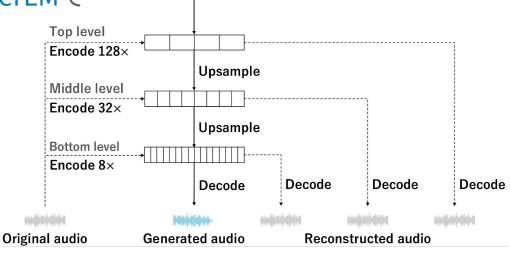
- Jukebox
- MusicGen

Jukebox: 階層的な音楽生成モデル [Dhariwal, 2020]

- 3つのレベルのVQ-VAE[Van Den Oord,2017]によって音楽信号を離散トークン化
 - Top:最も圧縮率が高い.大域な情報を捉える.
 - Middle:中間の圧縮率
 - Bottom:最も圧縮率が低い.局所的な情報を捉える.

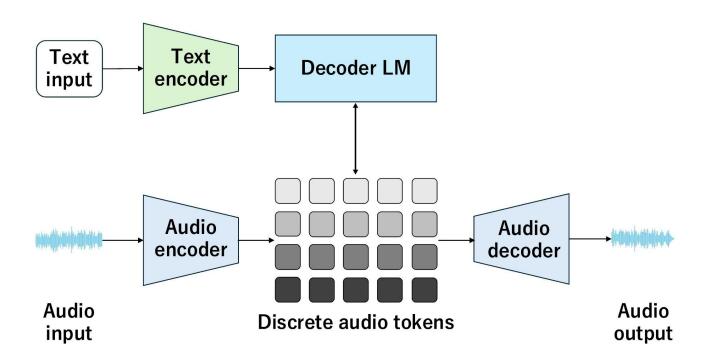
TransformerベースのDecoderLMで

階層的にトークン生成



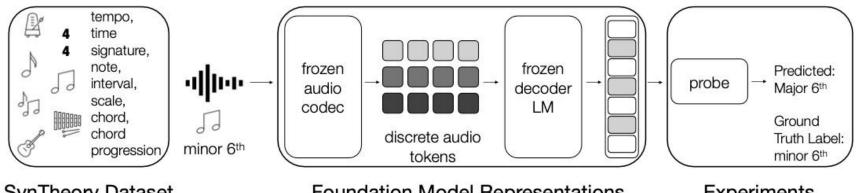
MusicGen:非階層的な音楽生成モデル [Copet, 2023]

- EnCodec [D´efossez, 2022]によって音楽信号を離散トークン化
- DecoderLMでトークン列生成



音楽の基本概念に対しての解析 [Wei, 2024]

- テンポや音高,コードなど音楽の基本概念を対象
- モデルから抽出した中間表現を小規模な識別器に入力して学習
- 識別器の精度を確かめ,モデルが「どれだけ基本概念情報を保持 しているか」を検証
- しかし、「どのような構造で表現されているか」は未知



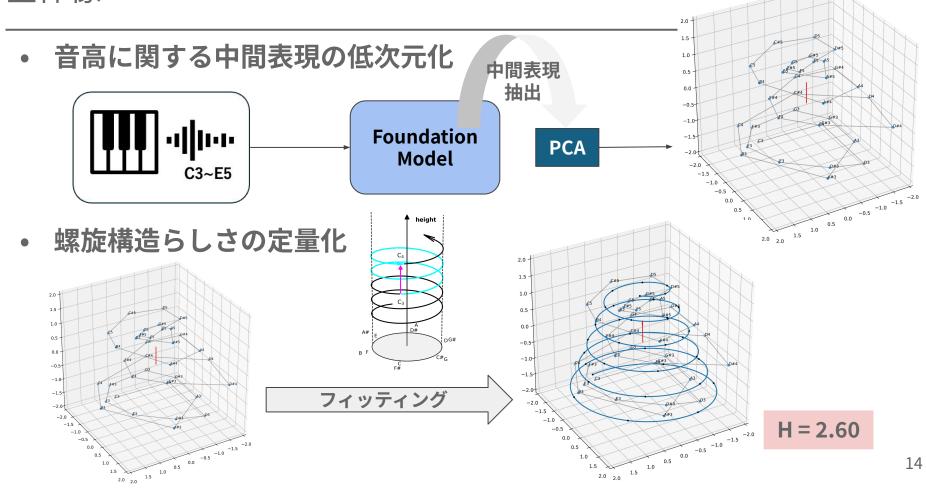
SynTheory Dataset

Foundation Model Representations

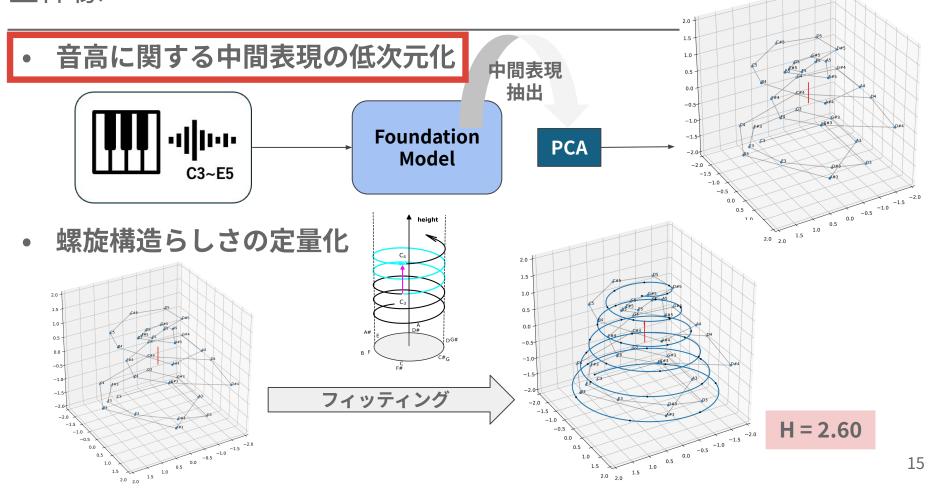
Experiments

提案手法

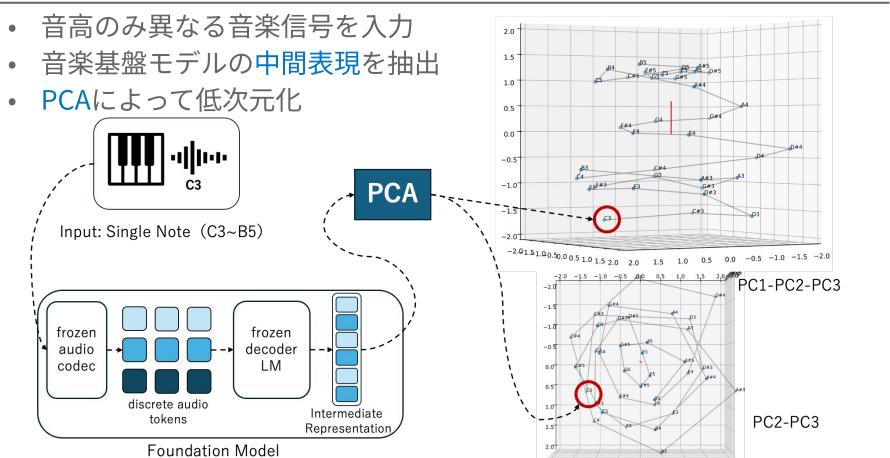
全体像



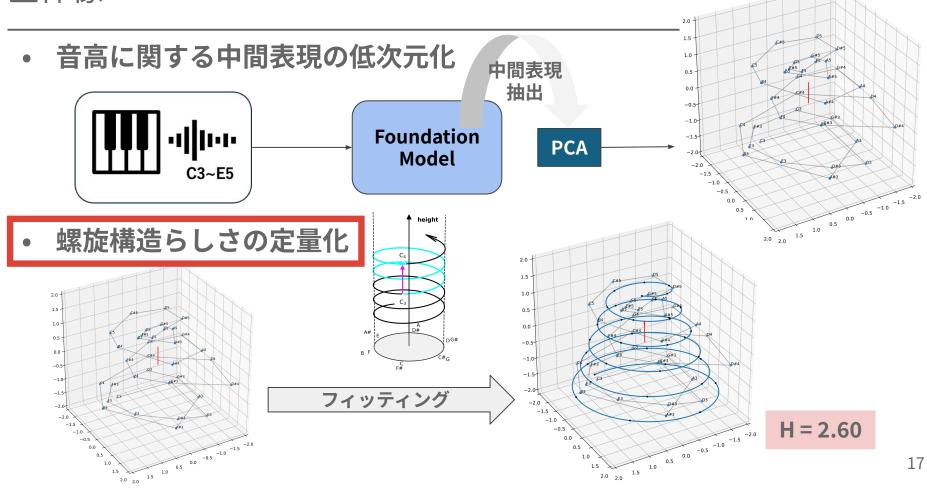
全体像



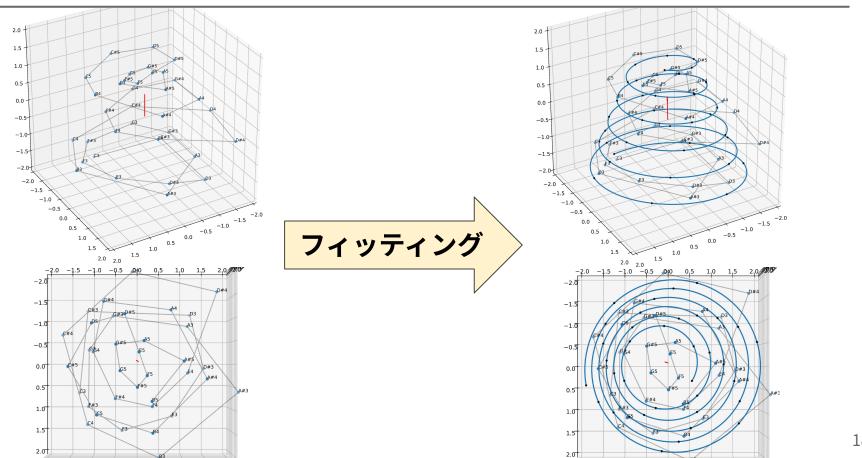
中間表現の獲得→次元削減 [Wei, 2024]



全体像



螺旋フィッティング例



螺旋フィッティング関数

高さ方向

回転方向

$$m{y}(t) = h(t) \cdot m{c} + r(t) \{\cos\theta(t) \cdot m{u} + \sin\theta(t) \cdot m{v}\}$$
 高さ 半径 位相 $\{m{c}, m{u}, m{v}\}$: 正規直交基底

t:音高インデックス

$$h(t) = h_{\text{pitch}} \cdot t + h_0$$

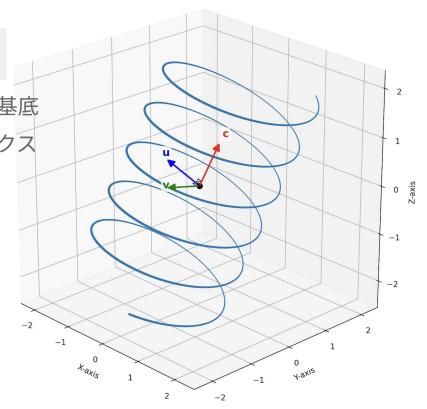
高さ変化係数 初期高さ

$$r(t) = r_{\text{slope}} \cdot t + r_0$$

半径变化係数 初期半径

$$\theta(t) = \omega_{\text{chroma}}(t - t_0)$$

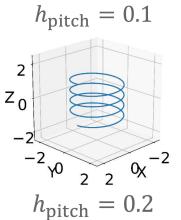
角周波数 回転の初期位置



パラメータ比較

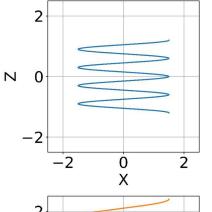
$$y(t) = h(t) \cdot c + r(t) \{\cos \theta(t) \cdot u + \sin \theta(t) \cdot v\}$$

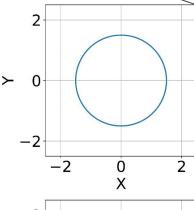
基準値 h_pitch = 0.15, ω _chroma = $\pi/3$, c:z軸, r0 = 1.5, r_slope = 0 h0 = 0, t0 = 0

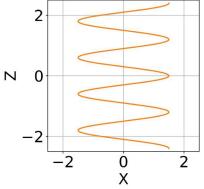


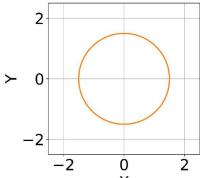
2

 z_0











$$r(t) = r_{\text{slope}} \cdot t + r_0$$

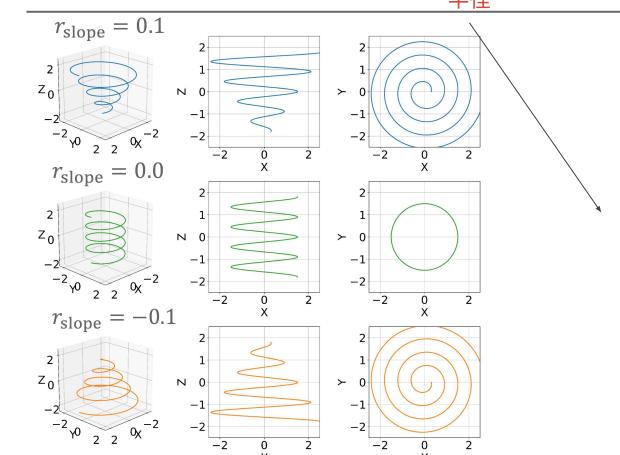
$$\theta(t) = \underline{\omega_{\text{chroma}}(t - \underline{t_0})}$$

角周波数 回転の初期位置

パラメータ比較

 $y(t) = h(t) \cdot c + r(t) \{\cos \theta(t) \cdot u + \sin \theta(t) \cdot v\}$

基準値 h_pitch = 0.15, ω_chroma = π/3, c:z軸, r0 = 1.5, r_slope = 0 h0 = 0. t0 = 0





$$r(t) = r_{\text{slope}} \cdot t + r_0$$
 半径変化係数 初期半径

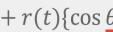
$$\theta(t) = \omega_{\rm chroma}(t - t_0)$$

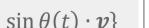
角周波数 回転の初期位置

パラメータ比較





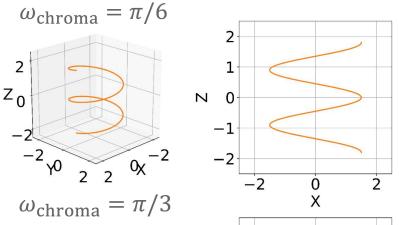




位相

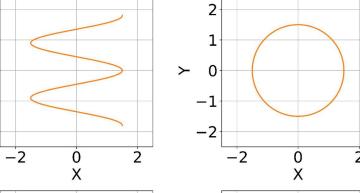


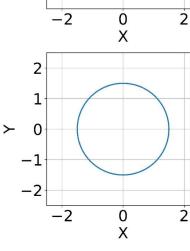
 ω _chroma = $\pi/3$, c:z軸. r0 = 1.5, $r_slope = 0$ h0 = 0, t0 = 0

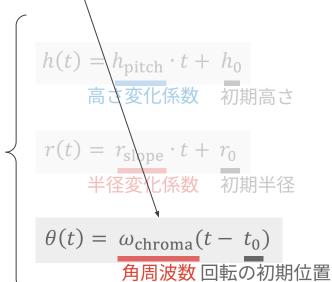


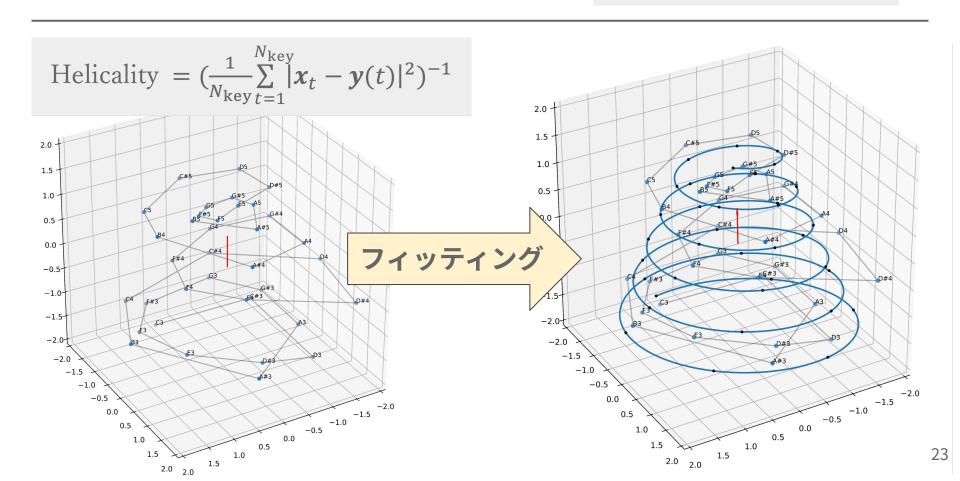
N 0

-2



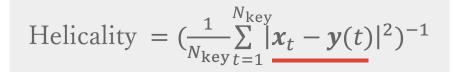




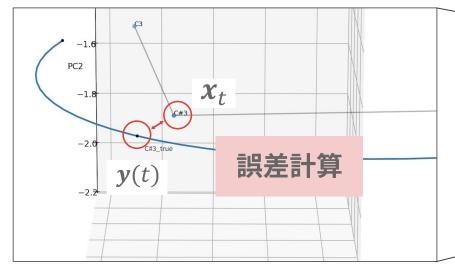


Helicalityスコア

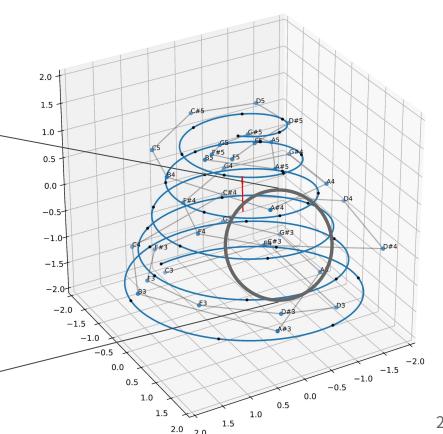
スコアが高い=綺麗な螺旋

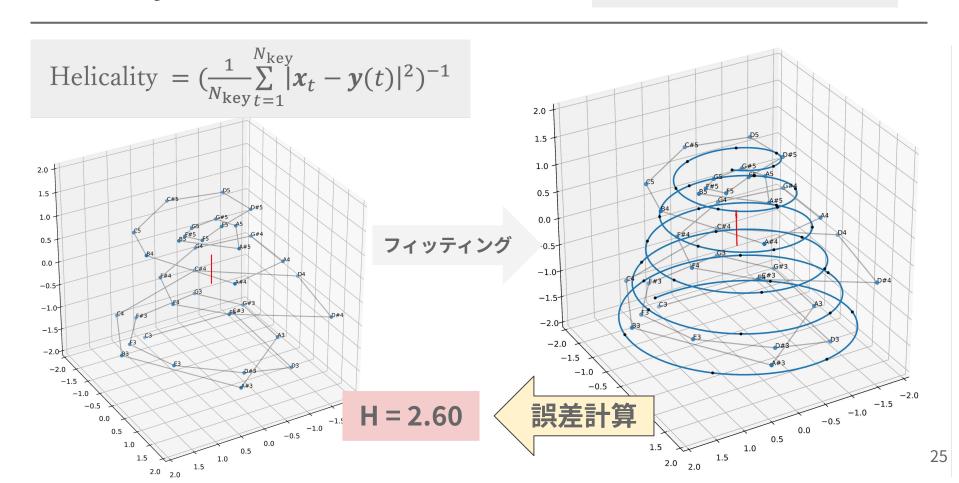


例) C#3の場合



$$y(t) = h(t) \cdot c + r(t) \{\cos \theta(t) \cdot u + \sin \theta(t) \cdot v\}$$



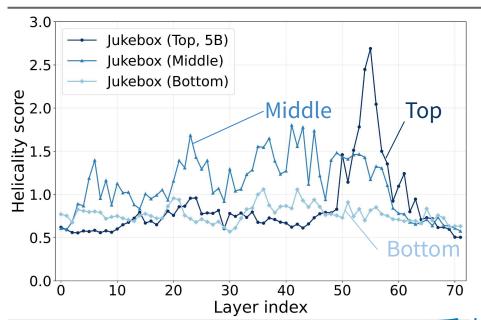


実験的評価

実験条件

モデル	Jukebox	Top(5B)	72層, 4800次元, 1トークン = 2.9ms	
		Top(1B)	72層, 4800次元, 1トークン = 2.9ms	
		Middle(1B)	72層, 1920次元, 1トークン = 0.73ms	
		Bottom(1B)	72層, 1920次元, 1トークン = 0.18ms	
	MusicGen	Small(300M)	24層, 1024次元, 1トークン = 20ms	
		Medium(1.5B)	48層, 1536次元, 1トークン = 20ms	
		Large(3.3B)	48層, 2048次元, 1トークン = 20ms	
データ	SynTheory [Wei, 2024]	Notes	C3~B5 までの3 オクターブ, ピアノの合成音源	
フィッティング	Optuna [Akiba, 2019]	試行回数1000 回, 繰り返し3回 → 最も良いHelicalityスコアを採用		

Jukeboxレベル別の比較



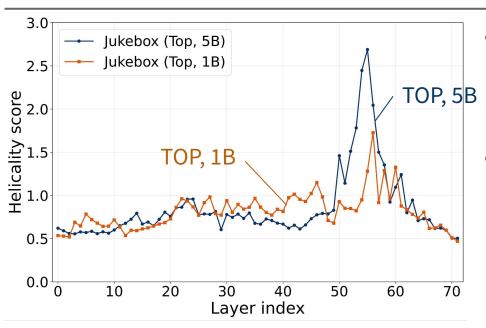
- 最大値はTopが持つ
 - Max: 2.6891 (Layer 55)
 - Mean: 0.8510
- 平均値はMiddleが最も高い
 - Max: 1.7138 (Layer 42)
 - Mean: 1.0378
- Bottomは全体的にスコアが低い
 - Max: 1.0584 (Layer 36)
 - Mean: 0.7714

- **Top:** $1 \vdash -0 = 2.9 \text{ms}$
- Middle: $1 \vdash 0 \lor = 0.73 \text{ms}$

Bottom: $1 \vdash - 0 = 0.18 \text{ms}$

- Topは受容野が広く,信号を粗く圧縮
 - →音高の階層構造を獲得
- Bottomは<mark>受容野が狭く</mark>,音楽波形の再現を担う
 - → 抽象的な階層構造を獲得しない

モデルサイズ別の比較(Jukebox)



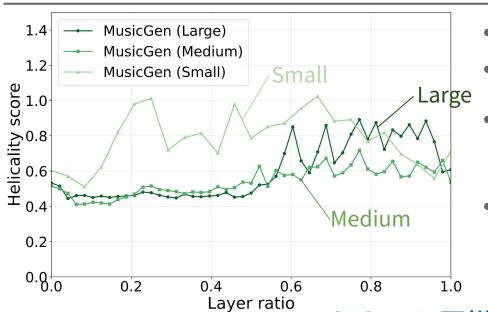
- 5Bモデル
 - 層の後半ほどスコアが上昇
 - 第55層付近で非常に高いスコア
- 1Bモデル
 - 分布は5Bと似ているが、ピークの スコアは控えめ

5B: Mean: 0.8510

1B: Mean: 0.8129

- モデルサイズが大きいほど、音高幾何に特化 する層が生まれる可能性
- Top層の後半で音高情報を埋め込んでいる

モデルサイズ別の比較 (MusicGen)



- Smallが全体的に高スコアを示す
- L/Mは全体的にスコアが低い
- モデルに関わらず中~後半層で ピークが出現する傾向
- モデルサイズが大きくなるほど、 ピークが明確に出現

Jukeboxと同様の傾向を持つ

- →「モデルサイズが大きくなるほど,音高幾何に 特化した層が生まれる」と言う仮説を支持
- Jukeboxとは対照的に, Smallが最も高いピーク

Large: 49層, 300M Medium: 49層, 1.5B

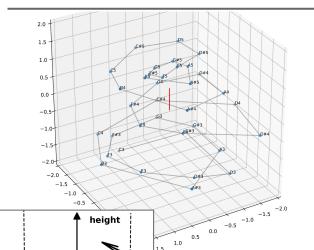
Medium: 49僧, 1.5b

Small: 25層, 3.3B

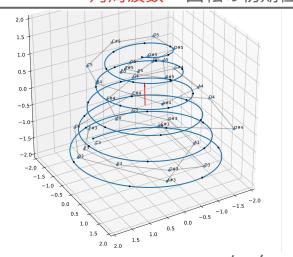
パラメータ分析:トライトーン



角周波数 回転の初期位置



フィッティング



- フィッティングパラメータ: $ω_chroma = 1.008 ≒ π/3 (1.047)$ → 1オクターブで**2回転**する螺旋のフィッティング
- 全音3つ分の関係にある音(トライトーン)は螺旋上で 円の対極に位置するのではなく,円の同じ位置となる →トライトーン代理などの役割を捉えている?

中心軸cとpitch heightの関係

- ・中心軸c へ射影した値と実際の 音高の順序との間の相関係数を算出
- Jukeboxは全ての階層において 高く安定した相関を示す
- Bottom階層ほど相関が高い
 - → 周波数を捉える能力は, 受容野に依存する可能性

受容野

Top: $1 \vdash - 0 = 2.9 \text{ms}$

Middle: $1 \vdash - 0 = 0.73 \text{ms}$

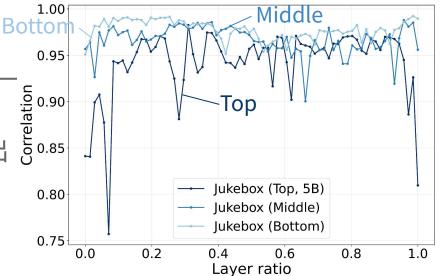
Bottom: $1 \vdash - 0 = 0.18 \text{ms}$

各層の相関係数の平均値

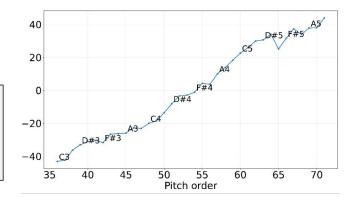
Top: Mean: 0.942

Middle: Mean: 0.965

Bottom: Mean: 0.977

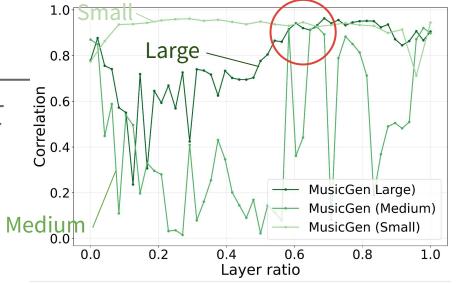


Jukebox bottomの第72層目 中心軸c への射影と音高の関係図(r = 0.992)



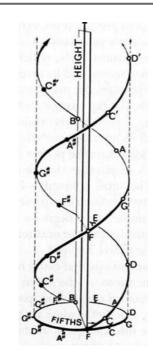
中心軸cとpitch heightの関係

- MusicGenはモデルサイズによって 振る舞いが異なる
- SmallはJukeboxに匹敵するほど 高く安定した相関を示す
 - → 限られたモデル容量の中で, 音楽の最も其木構造であるp
 - 音楽の最も基本構造であるpitch heightを優先的に学習した可能性
- Large/Mediumは平均相関が大幅に低い(L: 0.762, M: 0.440)
- 相関の最大値はSmallに匹敵する (S: 1.000, M: 0.996, L: 1.000)
 - → 「モデルサイズが大きくなるほど,音高幾何に特化した層が生まれる」 ということに対応する



今後の展望

- 別モデルの調査
 - MERTなど他の音楽基盤モデルとの比較
- その他構造の検証
 - 必ずしも螺旋構造になるわけではなさそう
 - 二重螺旋構造
 - 直線的構造
- 別の音色の検証
 - ピアノ以外の楽器音
 - 倍音を含まない純音
- テンポやコードなど他の要素への拡張
- 介入実験 → 制御性の検討



[Sheard, 1982]

まとめ

- 目的
 - 音楽基盤モデルが音高情報を螺旋構造に埋め込むかを調査する.
- 提案手法
 - 音高のみ異なる音楽データを入力とした時の中間表現をPCAで低次元化.
 - 螺旋モデルをフィッティングし,螺旋構造らしさを定量評価.
- 評価結果
 - 音高のような抽象的な概念はモデルの**深い層**でより明確に表現される.
 - モデルサイズが大きくなるほど、音高幾何に特化した層が生まれる.
- 結論
 - 音楽基盤モデルは**人間の音高知覚と類似した螺旋構造を内部に獲得**し、 モデルの**階層構造やサイズ**が螺旋表現の明瞭さと局所性に影響する.

補足

各モデルの最高スコア層におけるパラメータ一覧

モデル	層	Н	r0	h_pitch	r_slope	ω_chroma	h0	t0
Jukebox(Top)	55/72	2.69	1.77	0.095	-0.025	1.008	-1.588	0.238
Jukebox(Middle)	41/72	1.80	2.03	0.094	-0.050	-1.023	-1.545	-0.103
Jukebox(Bottom)	36/72	1.06	1.15	0.102	-0.006	1.058	-1.866	-4.313
MusicGen(Large)	37/48	0.89	1.40	0.088	-0.039	-0.567	-1.557	5.430
MusicGen(Medium)	37/48	0.72	1.42	0.061	-0.036	0.525	-1.102	-0.367
MusicGen(Small)	16/24	1.03	1.53	0.088	-0.026	-1.041	-1.482	3.998

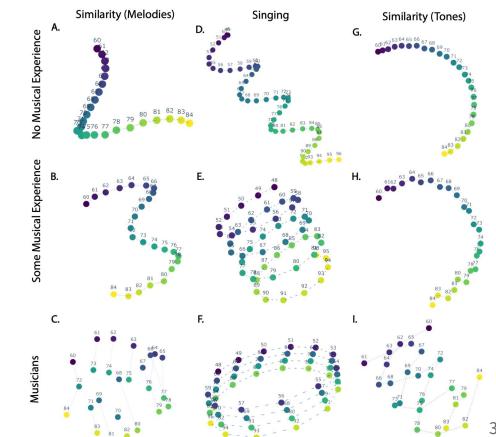
人間の音高知覚の研究

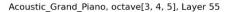
[Raja, 2023] Probing the Structure of Musical Pitch Across Tasks and Experience

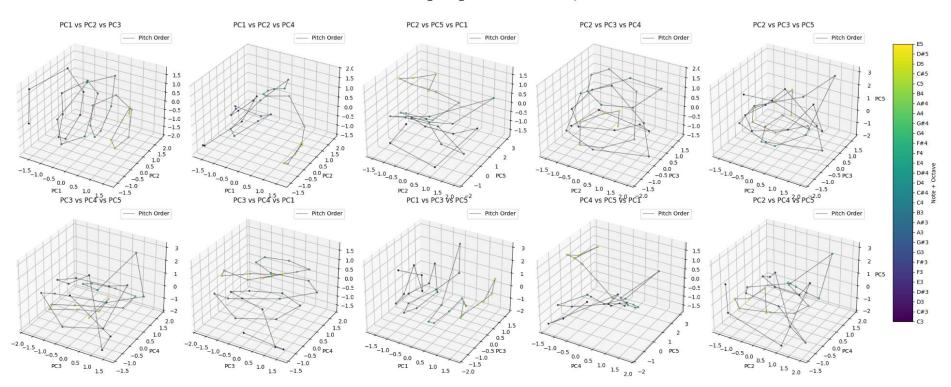
音楽経験によって現れる構造が異なる

熟練者ほど音楽的性質を 反映した形

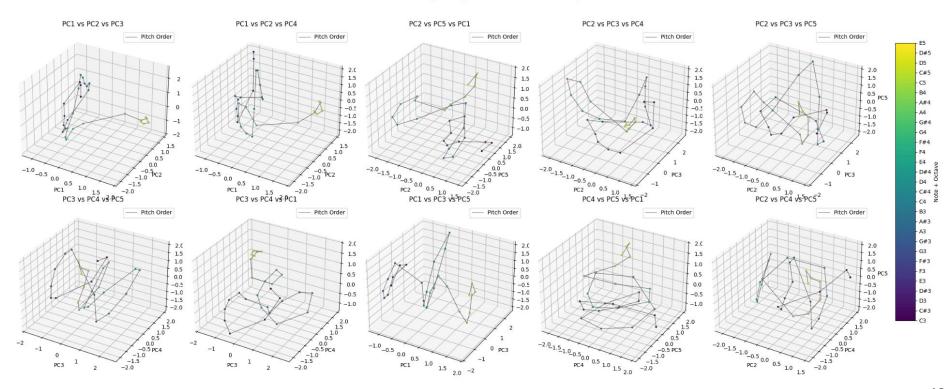
音楽家の螺旋構造の 高さ方向がほぼない → オクターブ識別力 が高い



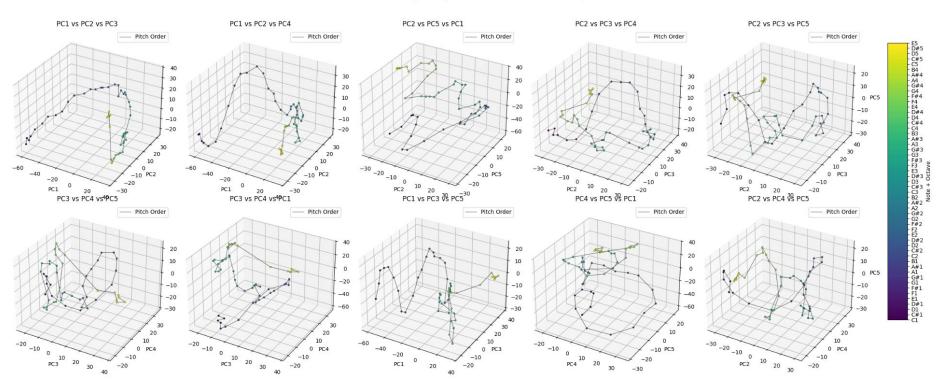




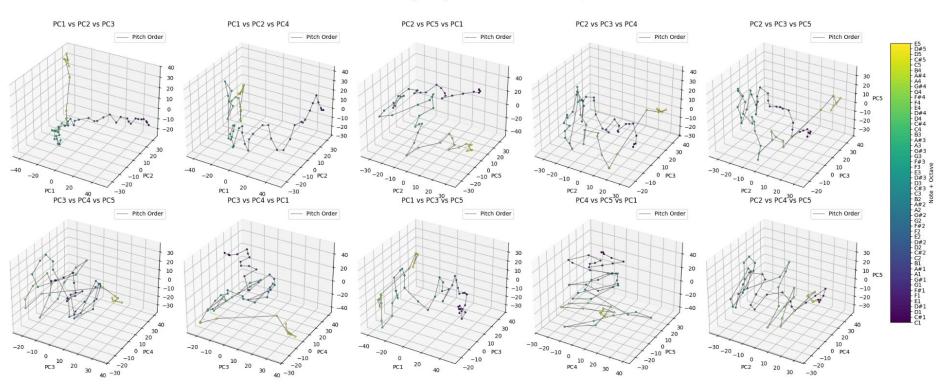
Acoustic_Grand_Piano, octave[3, 4, 5], Layer 48



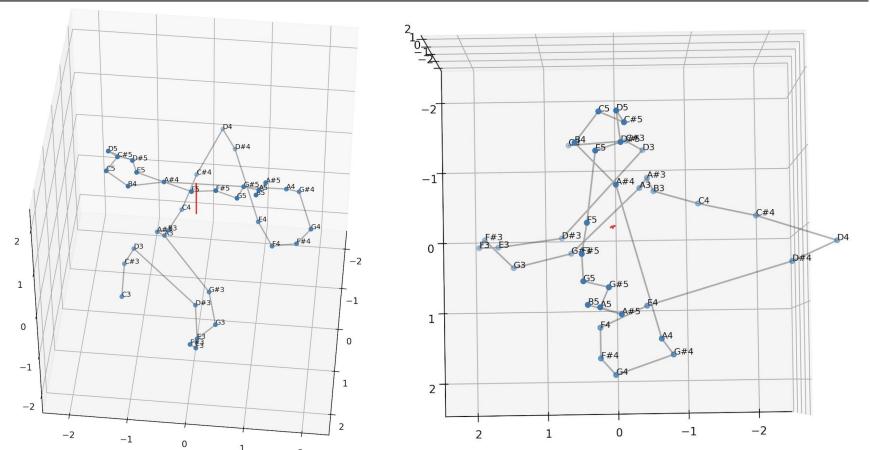




Acoustic_Grand_Piano, octave[1, 2, 3, 4, 5], Layer 36

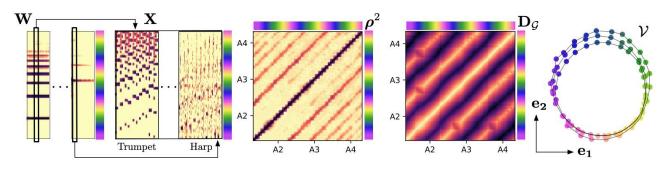


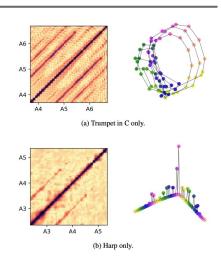
例3 Jukebox top 47層 pc1-pc2-pc3



音高の螺旋表現 (オクターブ等価性)

音響信号から音高螺旋構造を可視化 [Vincent, 2020]

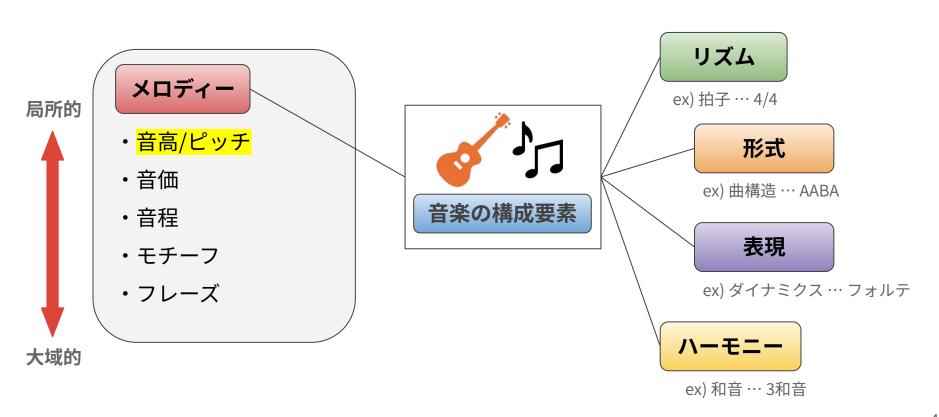




- 音響信号から音高螺旋構造を自動発見する手法
- 様々な楽器の単音データにこの手法を適応した結果,綺麗な螺旋を形成
- 音色によって構造が変わる
 - トランペットのような倍音が豊かな楽器は,綺麗な螺旋が形成される
 - ハープのような倍音が弱い楽器では、隣接音との相関が強く、直線的な構造

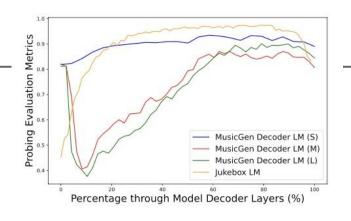
44

音楽の構成要素



音楽基盤モデルの解析

- 多くの音楽理論タスクで高性能
- モデルが音楽の基本概念を中間表現に エンコードしているとわかる



	Notes	Intervals	Scales	Chords	Chord Progressions	Tempos	Time Signatures	Average
Jukebox LM	0.951	0.995	0.978	0.997	0.971	0.993	1.000	0.984
MusicGen LM (S)	0.897	0.995	0.949	0.990	0.942	0.969	0.911	0.950
MusicGen LM (M)	0.851	0.983	0.863	0.989	0.870	0.956	0.883	0.914
MusicGen LM (L)	0.866	0.972	0.905	0.989	0.901	0.965	0.905	0.929
MusicGen Audio Codec	0.729	0.965	0.383	0.879	0.330	0.947	0.677	0.701
Mel Spectrogram	0.712	0.995	0.897	0.988	0.723	0.785	0.827	0.847
MFCC	0.467	0.822	0.370	0.863	0.872	0.923	0.688	0.715
Chroma	0.954	0.820	0.989	0.994	0.869	0.847	0.672	0.878
Aggregate Handcrafted	0.941	0.997	0.972	0.992	0.868	0.947	0.833	0.936