

1. はじめに

テキストからの2次元モーション生成の先行研究として、2CM-GPT[1] が提案されている。2CM-GPT は、3次元モーション生成モデルと異なり、収集が容易な2次元モーションのデータセットを学習に利用できる。そのため、生成に失敗した2次元モーションを収集して、データセットを動的に拡張することにより、効果的なファインチューニングが実現できる。一方で、2CM-GPTはいくつかの課題もある。1つ目は、各フレームのモーションを独立で生成するため、時間的な整合性が低い。2つ目は、テキストとモーションを対応付けることなく混在して学習するため、テキストとモーションの整合性が損なわれる。これらの課題を解決するために、本研究ではテキストとモーションをアテンション機構で関連付けるとともに、時間的な整合性を考慮する手法を提案する。

2. 2次元モーション生成

2CM-GPT は、2次元のモーションを生成する代表的な手法である。2CM-GPTのモデル構造を図1に示す。2CM-GPTは、Motion Tokenizerで人間のモーションを離散的なトークンに変換する。そして、テキストも同様にText Tokenizerでトークンに変換する。これらを連結させたMixed TokensをLanguage Encoderに入力して潜在ベクトルを獲得する。潜在ベクトルをもとにLanguage Decoderが出力したOutput TokensをMotion Tokenizerに入力して2次元のモーションを生成する。

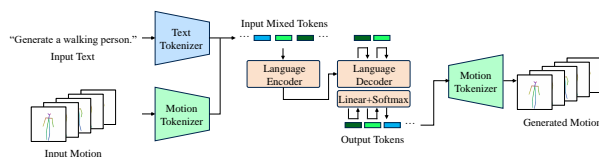


図1: 2CM-GPTのモデル構造

3. 提案手法

本研究では、2CM-GPTの学習アプローチが抱える「時間的な整合性」と「テキストとモーションの整合性」の不一致を解決する手法を提案する。提案手法のモデル構造を図2に示す。なお、提案手法がテキストと2次元モーションの関係性を学習する過程をTraining Phase、テキストから2次元モーションを生成する過程をT2M Phaseとする。

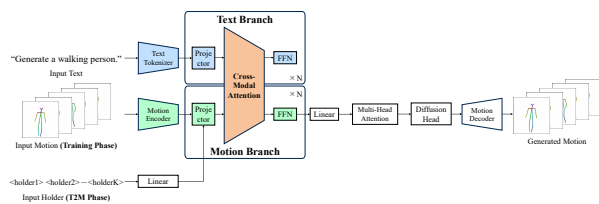


図2: 提案手法のモデル構造

3.1. 学習

Training Phaseでは、VAEで学習されたMotion Encoderを用いて人間の連続的な動作を離散化しない形で潜在空間に埋め込む。この埋め込み特徴をMotion Branchに入力する。テキストは、Text Tokenizerで埋め込み特徴とし、Text Branchに入力する。各ブランチにおいて、埋め込み特徴をProjectorに入力し、Query, Key, Valueベクトルを抽出する。そして、Motion BranchとText Branchの各ベクトルを同じCross-Modal Attentionに入力して、Motion BranchのValueベクトルにテキスト特徴を反映させる。Motion Branchの出力をMulti-Head Attentionに入力して潜在ベクトルを獲得する。潜在ベクトルをDiffusion Headに入力し、後述のT2M Phaseのための逆拡散過程を学習する。損失関数には、生成モーションと実モーションの差を評価する特徴再構成損失、対応関係にあるモー

ションとテキストがクロスモーダルな特徴空間上で近接するよう制約する分類損失、逆拡散過程における潜在表現の再構成誤差を評価する拡散損失を用いる。

3.2. モーション生成

T2M Phaseでは、Motion BranchとText BranchのCross-Modal Attentionによって、Motion Branchに入力したHolderにテキスト特徴を反映させる。Motion Branchの出力をMulti-Head Attentionに入力して、潜在ベクトルを獲得する。潜在ベクトルをDiffusion Headに入力して、逆拡散過程によるノイズ除去を行う。その後、VAEで学習されたMotion Decoderを用いて、潜在ベクトルから2次元のモーションを生成する。

4. 評価実験

2CM-GPTとの比較実験により、提案手法の有効性を示す。

4.1. 定量的評価

表1より、提案手法はFIDが低いことから、2CM-GPTと比べて2次元モーション生成精度の向上を確認した。一方、2CM-GPTのDiversityは提案手法よりも高い。その要因として、モーション生成精度が十分でないために、類似の指示文に対しても多様なモーションを生成することが考えられる。

表1: テキストからの2次元モーション生成精度の比較

Method	FID ↓		Diversity ↑	
	real	gen	real	gen
2CM-GPT	-1.37×10^{-9}	32.36	16.96	19.28
提案手法	-5.24×10^{-9}	12.15	16.69	12.92

そこで、多様性の要因を検証するために、表1のDiversityと、後述の定性的評価で用いる図3の指示文のみを与えて生成されたモーションのDiversityを比較する。表2より、2CM-GPTは同一の指示文のみを与えた場合でもDiversityが高い傾向を示したことから、上記の可能性を裏付ける傾向が確認された。

表2: 指示文ごとの生成モーションのDiversityの比較

Method	Multi Instruction	Single Instruction
2CM-GPT	19.28	18.85
提案手法	12.92	6.07

4.2. 定性的評価

2CM-GPTと提案手法に同じ指示文を与えて生成させた2次元モーションを図3に示す。図3の手の動きに注目すると、提案手法の生成モーションは2CM-GPTと比較して、指示文の動作内容を正確に反映していると判断できる。

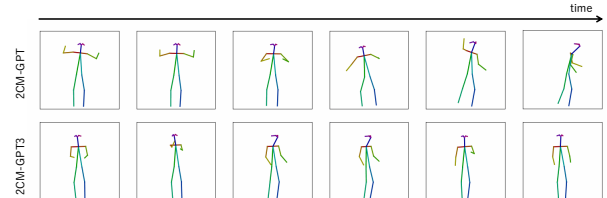


図3: テキストからの2次元モーション生成結果の可視化

5. おわりに

本研究では、2CM-GPTと提案手法の評価実験を行い、提案手法の有効性を示した。今後は、提案手法で生成したモーションを用いてポーズ誘導による人物動画生成を実施し、実用性を検証する。

参考文献

- [1] R. Inoue, *et al.*, “2D Motion Generation Using Joint Spatial Information with 2CM-GPT”, VISIGRAPP, vol.2, pp.582-590, 2025.