

When Counterpoint Meets Chinese Folk Melodies

Nan Jiang[†], Sheng Jin[†], Zhiyao Duan[‡], Changshui Zhang[†]

[†] Institute for Artificial Intelligence, Tsinghua University (THUI),
State Key Lab of Intelligent Technologies and Systems,
Beijing National Research Center for Information Science and Technology (BNRist),
Department of Automation, Tsinghua University, Beijing, China
[‡] Department of Electrical and Computer Engineering, University of Rochester



Acknowledgement:

National Key R&D Program of China (No. 2018AAA0100701),
Beijing Academy of Artificial Intelligence (BAAI),
National Science Foundation grant No. 1846184.

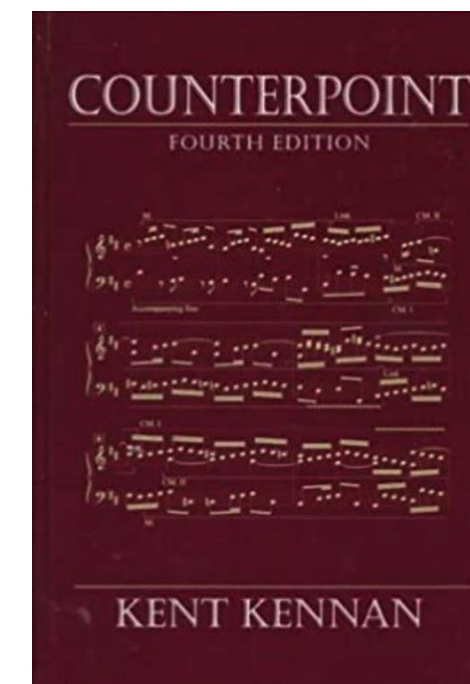
Introduction

Human-machine collaborative duet improvisation.

- ✓ Chinese folk melody style
- ✓ Counterpoint interaction between parts



Task: Incorporating Western counterpoint interactions into Chinese folk melodies for online human-machine collaborative duet improvisation.



- **Chinese folk melody:** typically presented in a *monophonic* form or with accompaniments that are less melodic.
- **Counterpoint:** mediation of two or more musical voices into a meaningful and pleasing whole.

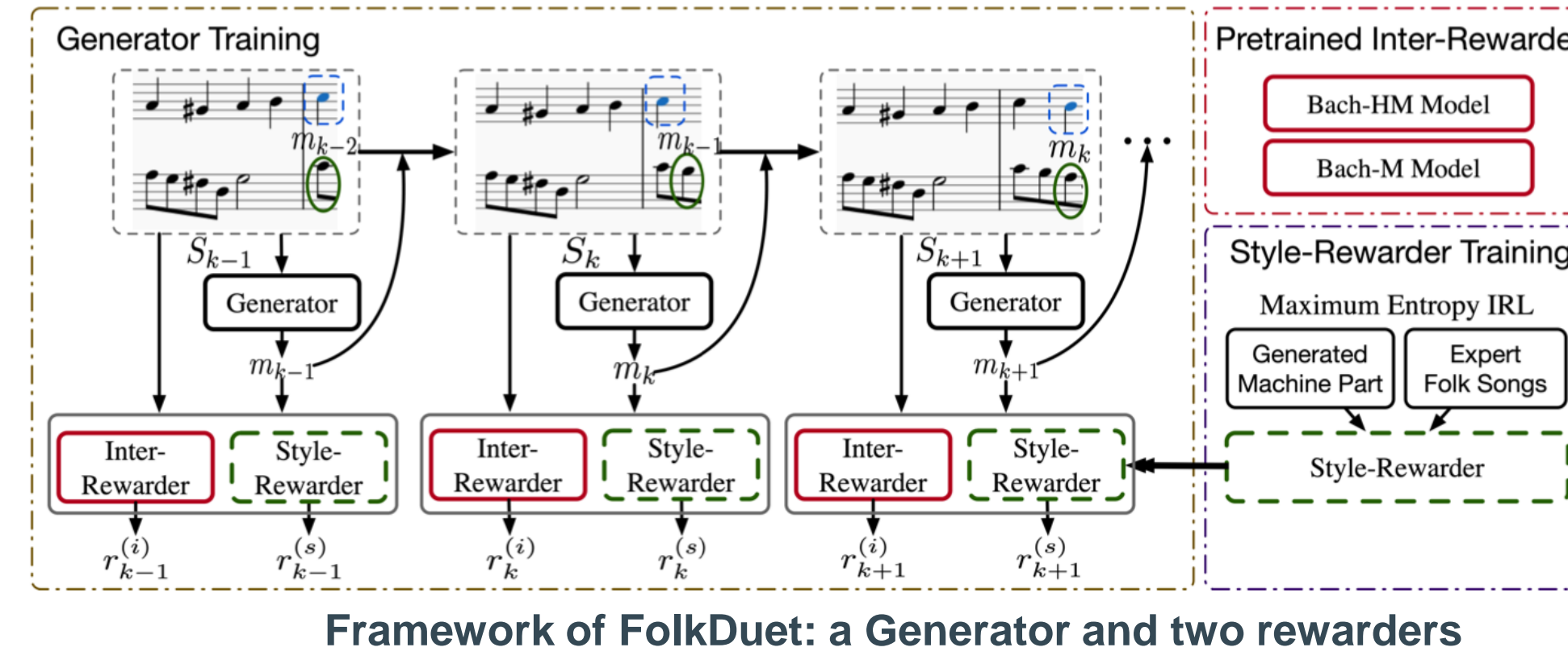
Challenges

- Out-of-domain data (Chinese folk duets are scarce)
Monophonic Chinese folk melodies + Bach chorales
- Counterpoint pattern is coupled with western-music style patterns
Extract counterpoint interaction pattern & eliminate Western music style

Our Solutions

- Reinforcement Learning → Design task-specific reward functions
- Measure counterpoint interaction using *mutual information*

FolkDuet



Inter-Rewarder: models the counterpoint interaction in Western music
Style-Rewarder: models the melodic pattern of Chinese folk melodies
The Generator is trained using reinforcement learning with these two rewards.

Rewarders

Inter-Rewarder: measures the degree of interaction between human and machine parts through a mutual information informed measure.

$$I(X, Y) = \sum_{X, Y} P(X, Y) \log \frac{P(X, Y)}{P(X)P(Y)}$$

$$= \sum_{X, Y} P(X, Y) [\log P(Y|X) - \log P(Y)] \approx \sum_{X_i, Y_i \sim P_{X, Y}} [\log P(Y_i|X_i) - \log P(Y_i)]$$

$$= \sum_{X_i, Y_i \sim P_{X, Y}} \left[\log \prod_{k=1}^{K^y} P(y_k^{(i)} | X_i, y_{1:k-1}^{(i)}) \cdot P(y_k^{(i)} | X_i) - \log \prod_{k=1}^{K^y} P(y_k^{(i)} | y_{1:k-1}^{(i)}) \cdot P(y_k^{(i)}) \right]$$

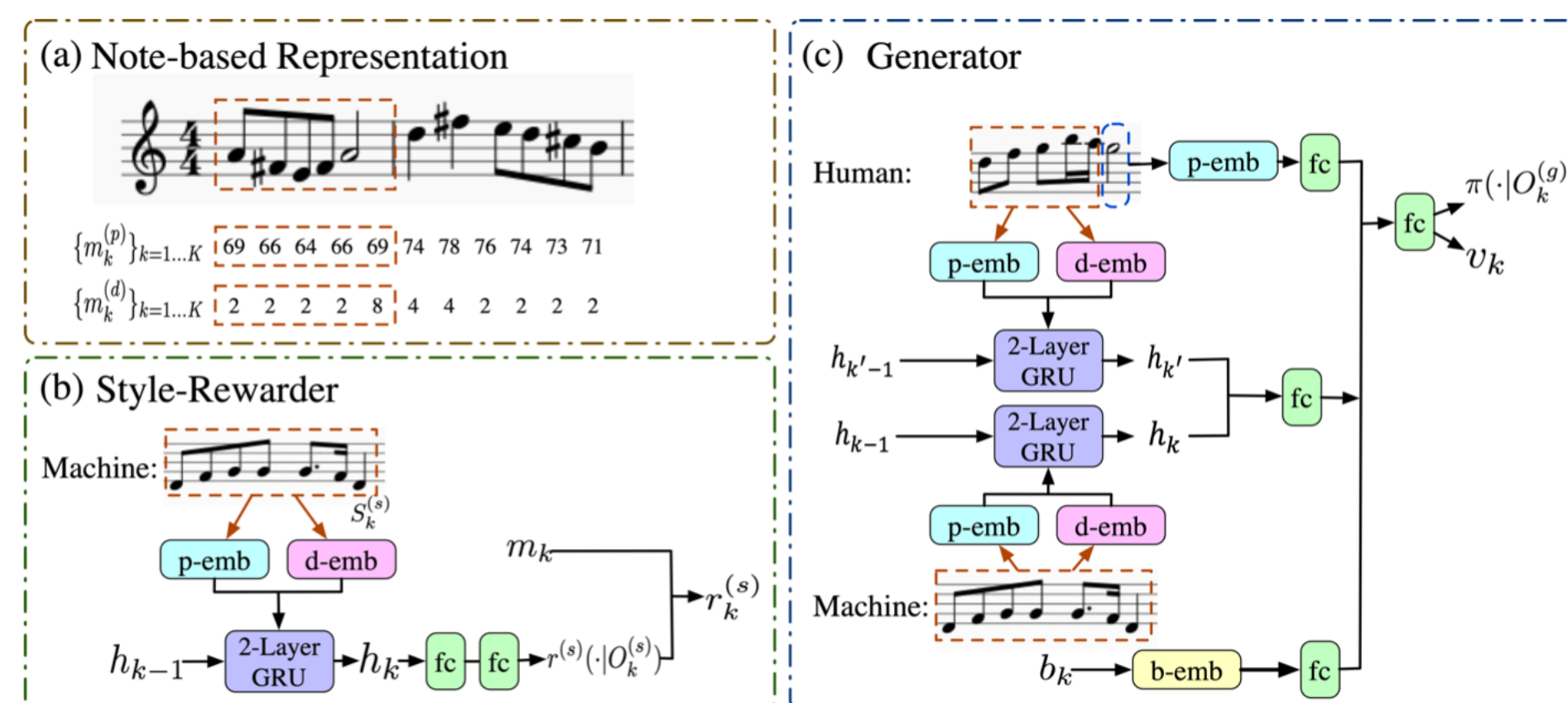
$$= \sum_{X_i, Y_i \sim P_{X, Y}} \sum_{k=1}^{K^y} [\log P(y_k^{(i)} | X_i, y_{1:k-1}^{(i)}) - \log P(y_k^{(i)} | y_{1:k-1}^{(i)})] + C(X_i, y_k^{(i)})$$

$$\log p(\text{Machine}|\text{Human}) - \log p(\text{Machine})$$

Style-Rewarder: Inverse Reinforcement learning (IRL): learns to infer a reward function underlying the observed expert behavior.

Style-Rewarder is alternatively updated using the maximum entropy inverse reinforcement learning. Its learning objective is to infer the reward function that underlies the demonstrated expert behavior, *i.e.* the Chinese folk melodies.

Architectures



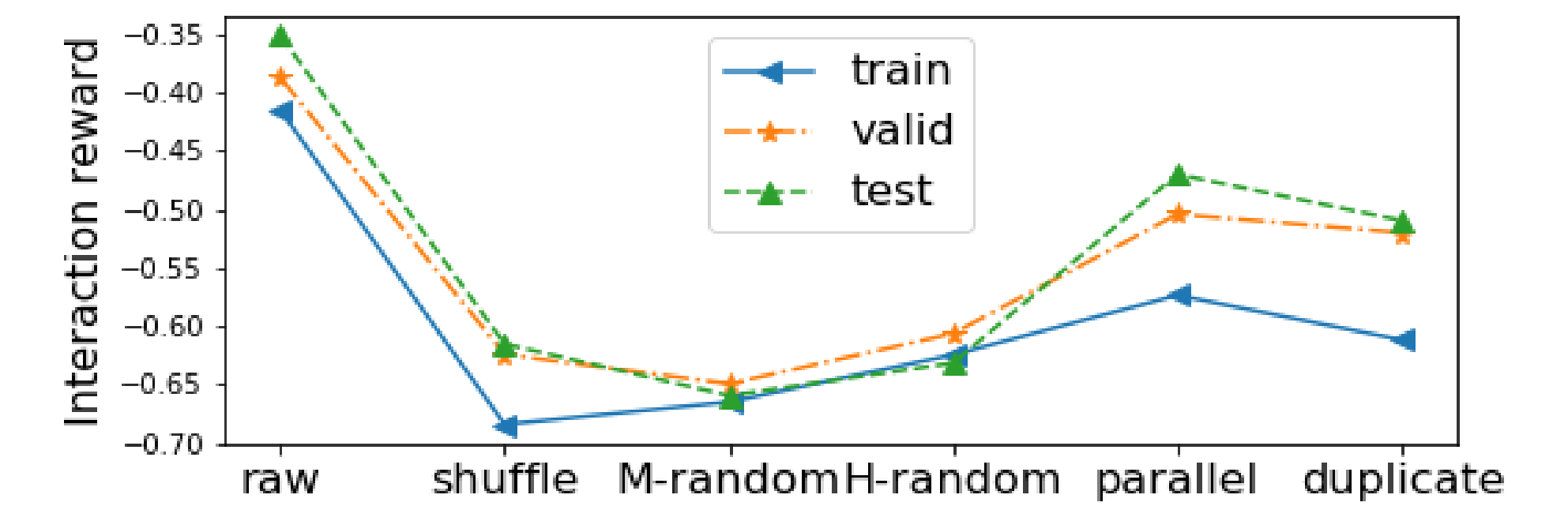
(a) The note-based representation, the network architectures of (b) Style-Rewarder and (c) Generator. p-emb, d-emb and b-emb represent pitch/duration/beat embedding modules, respectively. GRU represents the Gate Recurrent Unit, and fc stands for the fully-connected layer.

Results

Generated Duets



Can interaction reward reflect counterpoint interaction?



original **two randomly** **random** **random** **parallel H** **duplicate**
bach's duets **shuffled parts** **notes in** **notes in** **and M** **H and M**
the M part **the H part** **parts** **parts**

H part and M part are short for human part and machine part respectively. It shows that this interaction reward achieves the highest score on the original duets.

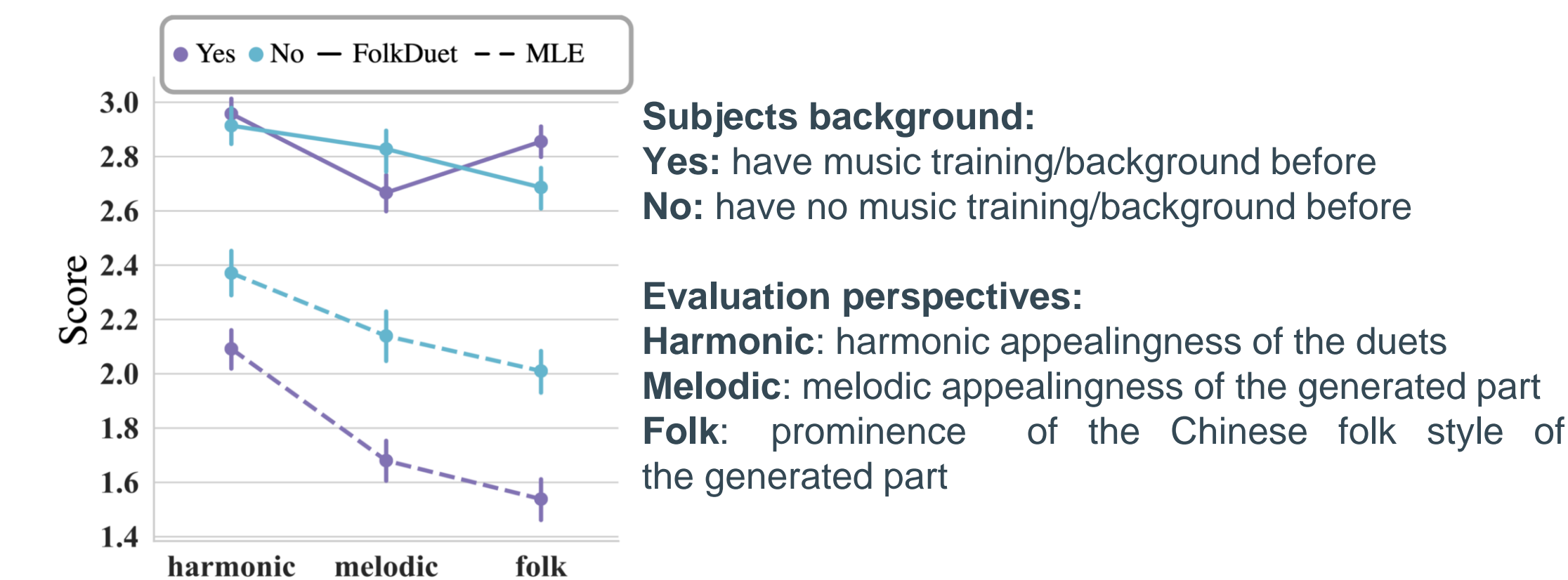
Objective Evaluation

Dataset	PC/bar	PI	IOI	PCH ↓	NLH ↓	key-consist ↑	inter-reward ↑
MLE	4.21 ± 0.12	3.02 ± 0.12	2.87 ± 0.10	0.017 ± 0.002	0.036 ± 0.008	0.78 ± 0.01	-0.30 ± 0.02
RL-Duet [27]	3.23 ± 0.01	4.02 ± 0.01	3.64 ± 0.02	0.017 ± 0.001	0.055 ± 0.002	0.71 ± 0.01	-0.50 ± 0.004
FolkDuet	3.96 ± 0.12	2.44 ± 0.14	2.16 ± 0.10	0.008 ± 0.001	0.014 ± 0.004	0.85 ± 0.01	0.13 ± 0.03

Style: Closer to Chinese folk datasets, in some statistics and distribution distance, e.g. pitch interval (PI), pitch class histogram (PCH).

Counterpoint interaction: Higher key consistency between human and machine parts, higher inter-reward.

Subjective Evaluation



Subjects background:
Yes: have music training/background before
No: have no music training/background before

Evaluation perspectives:
Harmonic: harmonic appealingness of the duets
Melodic: melodic appealingness of the generated part
Folk: prominence of the Chinese folk style of the generated part