# Sketching the Expression: Flexible Rendering of Expressive Piano Performance with Self-Supervised Learning

Seungyeon Rhyu[1], Sarah Kim[2], and Kyogu Lee[1]

[1]*Music and Audio Research Group, Seoul National University, South Korea*
[2]*Krust Universe, South Korea*

ISMIR 2022 BENGALURU

MARG Music & Audio Research Group

## Introduction

**Piano performance generation: Generating parameters that fit a musical piece**
- It aims to generate coherent parameters for loudness or timing, based on given musical scores.

**Expanding musical creativity: Musical expression beyond the written guidelines**
- Previous models for controlling a piano performance followed written expression guidelines or dealt with only partial attributes of musical expression.
- Performers can actively choose techniques to highlight various emotions or nuances, creating musical expressions beyond the guidelines.
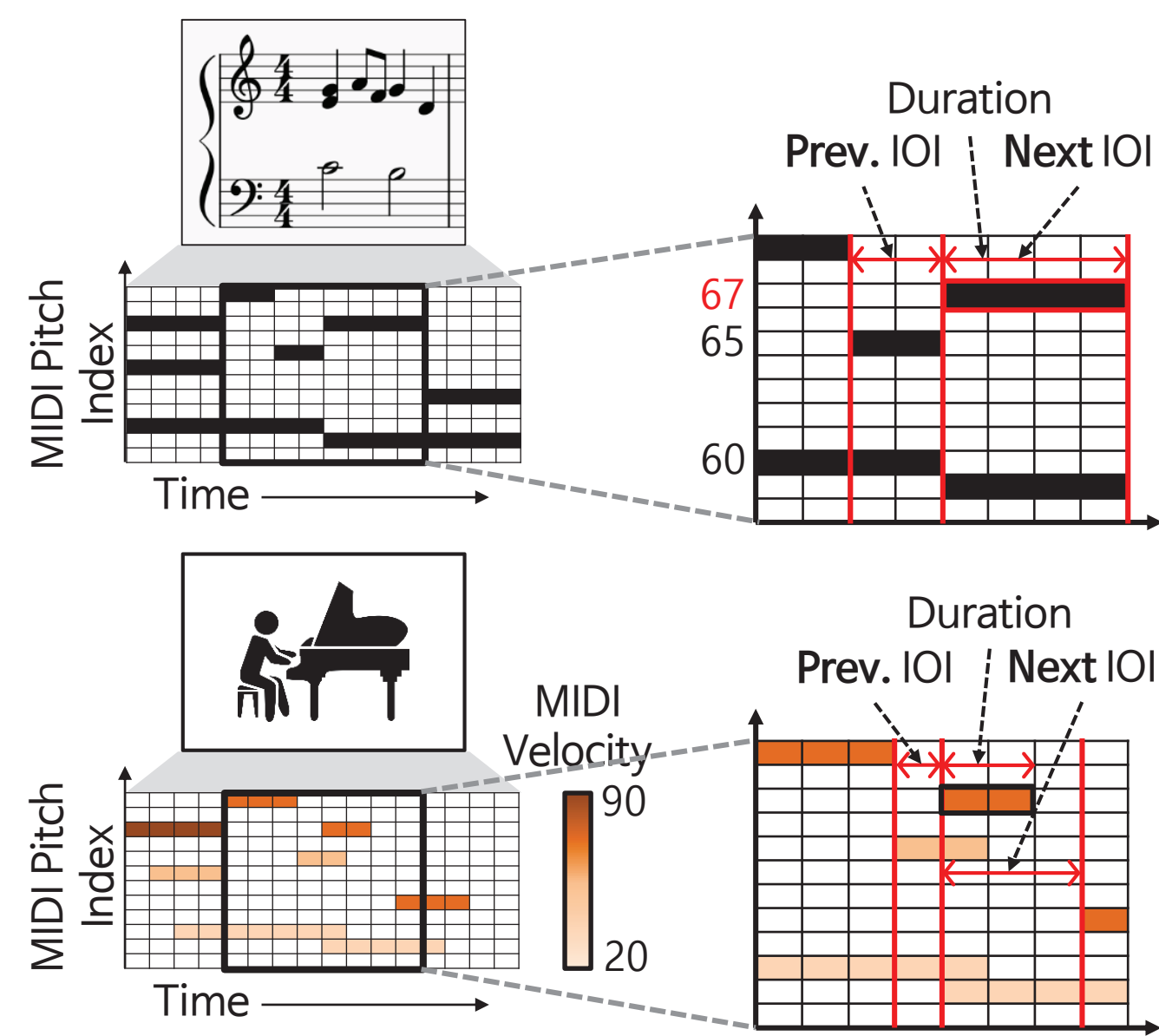
**Objective: Disentangling musical expression using self-supervised learning**
- We propose a generative model that disentangles two representations for high-level musical expression, or *explicit planning*, and low-level structural attributes.
- We use a conditional VAE modified for sequential data and a self-supervised learning framework to regularize the representations.
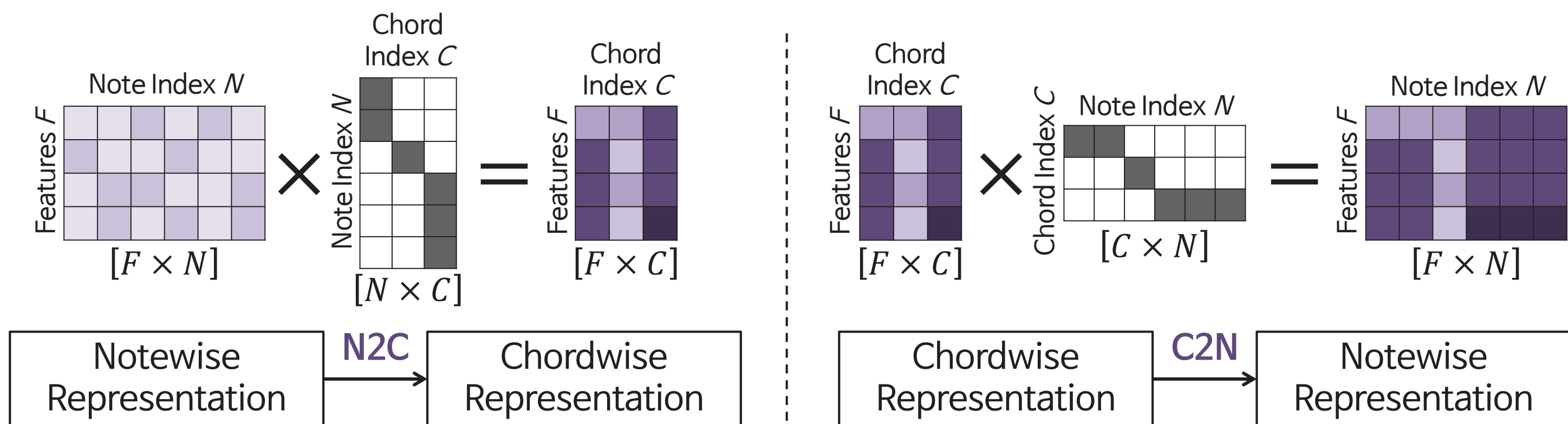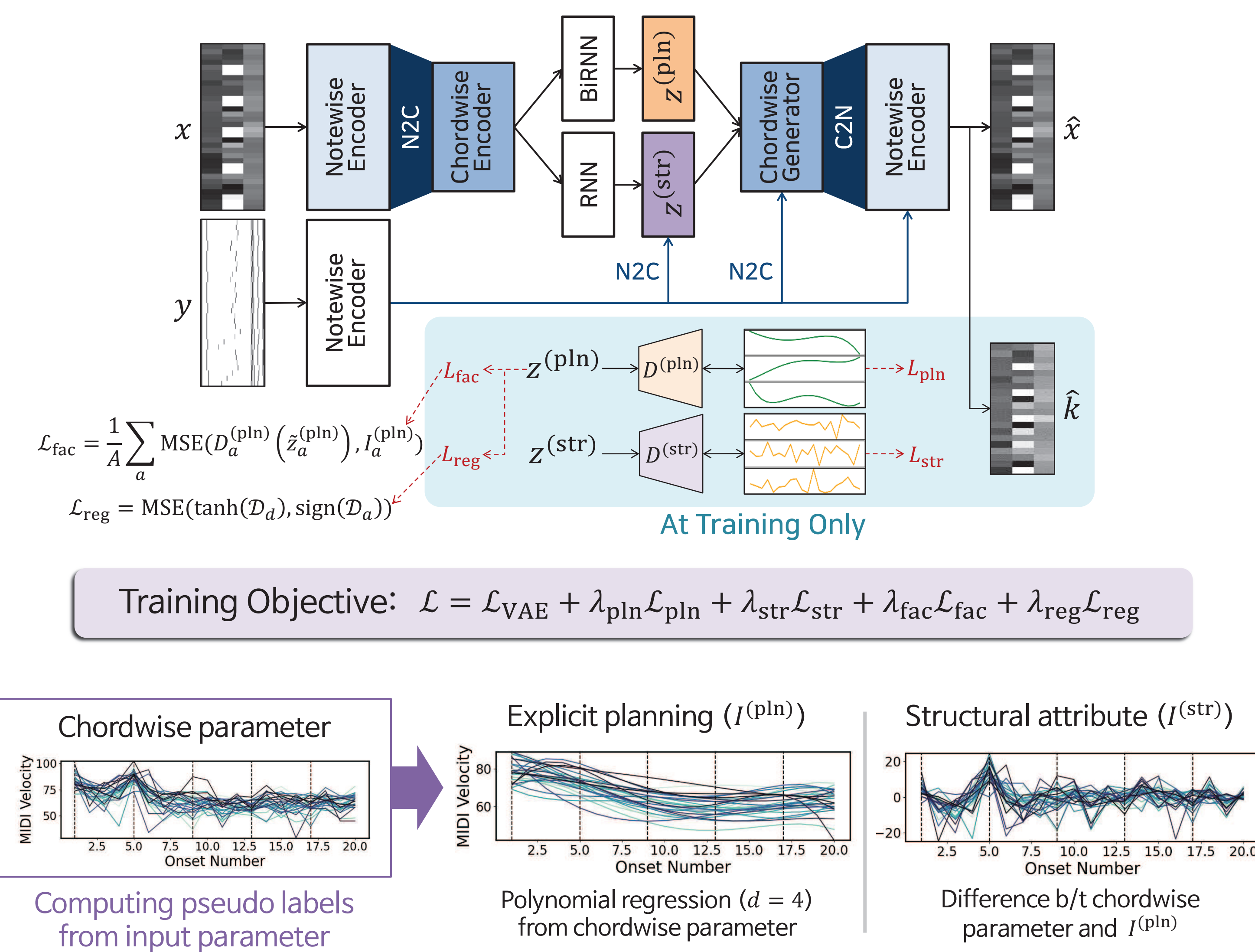
## Proposed Methods

### I. Data Representation

| Score Feature | Performance Feature |
|---|---|
| ✓ MIDI Pitch | ✓ MIDI Velocity (dynamics) |
| ✓ Duration (16th note) | ✓ IOI Ratio (tempo) |
| ✓ IOI (16th note) | $\frac{\text{Perform. Prev. IOI}}{\text{Score Prev. IOI}}$ |
| ✓ Is-top-voice | |
| ✓ #Note-in-Chord | ✓ Articulation |
| ✓ Position-in-chord | $\frac{\left(\frac{\text{Perform. Duration}}{\text{score Duration}}\right)}{\text{IOI Ratio with Next IOI}}$ |
| ✓ Staff | |
| ✓ Is-downbeat | |



### II. Modeling Musical Hierarchy



| Notewise Representation | N2C | Chordwise Representation |
|---|---|---|

| Chordwise Representation | C2N | Notewise Representation |
|---|---|---|

### III. Model Architecture



$$\mathcal{L}_{\text{fac}} = \frac{1}{A} \sum_a \text{MSE}(D_a^{(\text{pln})}(\tilde{z}_a^{(\text{pln})}), I_a^{(\text{pln})})$$

$$\mathcal{L}_{\text{reg}} = \text{MSE}(\tanh(\mathcal{D}_d), \text{sign}(\mathcal{D}_d))$$

At Training Only

Training Objective: $\mathcal{L} = \mathcal{L}_{\text{VAE}} + \lambda_{\text{pln}}\mathcal{L}_{\text{pln}} + \lambda_{\text{str}}\mathcal{L}_{\text{str}} + \lambda_{\text{fac}}\mathcal{L}_{\text{fac}} + \lambda_{\text{reg}}\mathcal{L}_{\text{reg}}$



Chordwise parameter → Explicit planning ($I^{(\text{pln})}$) → Structural attribute ($I^{(\text{str})}$)

Computing pseudo labels from input parameter | Polynomial regression ($d = 4$) from chordwise parameter | Difference b/t chordwise parameter and $I^{(\text{pln})}$

## Baseline Methods

### I. Architecture

| | Description |
|---|---|
| Notewise | Ours without chordwise encoding and decoding |
| CVAE | Variant of Notewise where $z^{(\text{pln})}$ is substituted with the supervisory signal $I^{(\text{pln})}$ |

### II. Ablation Study

| | Description |
|---|---|
| w/ $\mathcal{L}_{\text{pln}}$ | Ours only with prediction task using $z^{(\text{pln})}$ |
| w/ $\mathcal{L}_{\text{pln}} + \mathcal{L}_{\text{str}}$ | Ours with prediction tasks using $z^{(\text{pln})}$ and $z^{(\text{str})}$ |
| w/o $\mathcal{L}_{\text{fac}}$ | Ours without the additional factorization loss |
| w/o $\mathcal{L}_{\text{reg}}$ | Ours without regularization method* for sketch-control |

*A. Pati and A. Lerch. 2019. Latent space regularization for explicit control of musical attributes. *In Proceedings of the 36th International Conference on Machine Learning.*

## Dataset

| Evaluation | Objective | | Subjective |
|---|---|---|---|
| Type | Internal | External | External |
| Dataset | Yamaha e-Competition Vienna 4x22 Piano Corpus | ASAP | Online |
| Composer / Genre | Chopin only / Classical | 10 composers / Classical | Various / Non-Classical |
| # of song/perform. | 34 / 356 | 23 / 116 | 42 / (score only) |

## Evaluation

### I. Generation Quality

| Dataset | Internal | | | External | | |
|---|---|---|---|---|---|---|
| Metric | $R_{\text{recon}}$ | $R_{x\|\text{pln}}$ | $R_{x\|\text{pln}_0}$ | $R_{\text{recon}}$ | $R_{x\|\text{pln}}$ | $R_{x\|\text{pln}_0}$ |
| Notewise | 0.870 | 0.392 | 0.203 | 0.875 | 0.479 | 0.177 |
| CVAE | 0.730 | 0.338 | 0.223 | 0.741 | 0.399 | 0.216 |
| $L_{\text{pln}}$ | 0.627 | 0.357 | 0.229 | 0.687 | 0.414 | 0.220 |
| $L_{\text{pln}} + L_{\text{str}}$ | 0.770 | 0.325 | 0.181 | 0.837 | 0.398 | 0.195 |
| w/o $L_{\text{fac}}$ | 0.774 | 0.289 | 0.176 | 0.838 | 0.354 | 0.173 |
| w/o $L_{\text{reg}}$ | 0.737 | 0.437 | 0.224 | 0.793 | 0.502 | 0.216 |
| Ours | 0.737 | 0.427 | 0.231 | 0.789 | 0.498 | 0.203 |

- Pearson's correlation coefficients
- $R_{\text{recon}}$: Reconstruction loss.
- $R_{x\|\text{pln}}$: Evaluating samples with random $z_S$ and inferred $z^{(\text{pln})}$ from data.
- $R_{x\|\text{pln}_0}$: Evaluating samples with random $z_S$ and inferred $z_0^{(\text{pln})}$ from a zero matrix.

- Proposed architecture outperforms CVAE in most metrics.
- Proposed chordwise model generates better results with random $z^{(\text{str})}$ than Notewise.
- Ours shows stable generation scores with random $z^{(\text{str})}$ compared to other models.

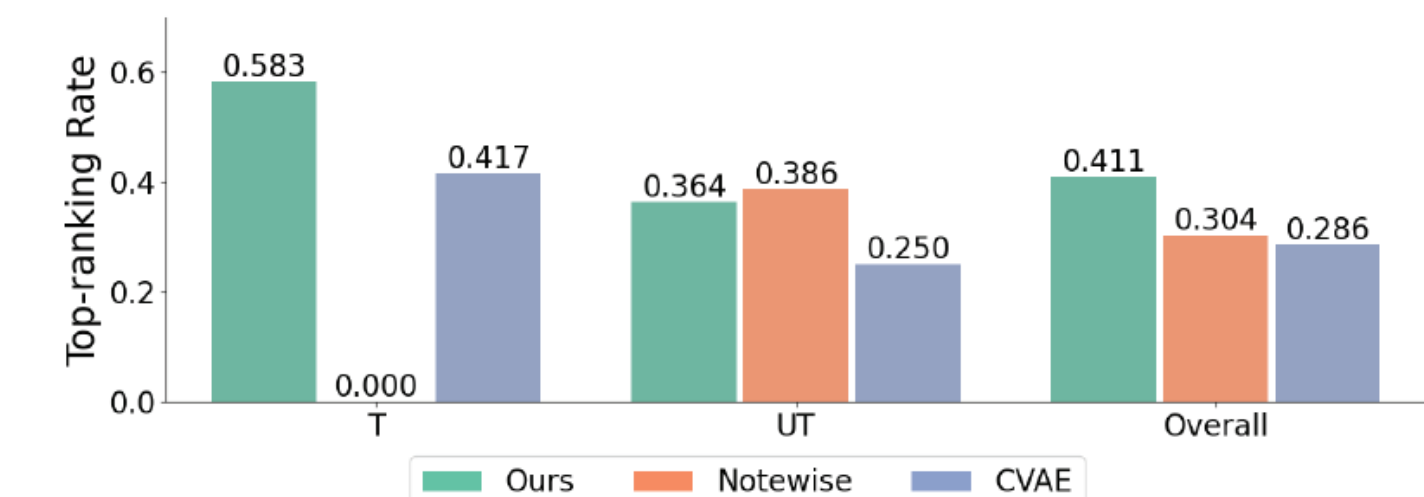### II. Disentanglement & Controllability of Musical Expression

| Dataset | Internal | | External | |
|---|---|---|---|---|
| Metric | $\text{MSE}_P$ | $\text{MSE}_S$ | $\text{MSE}_P$ | $\text{MSE}_S$ |
| Notewise | 0.003 | 0.006 | 0.022 | 0.028 |
| CVAE | 0.034 | 0.045 | 0.085 | 0.092 |
| $L_{\text{pln}}$ | 0.028 | 0.036 | 0.074 | 0.077 |
| $L_{\text{pln}} + L_{\text{str}}$ | 0.012 | 0.015 | 0.022 | 0.027 |
| w/o $L_{\text{fac}}$ | 0.018 | 0.023 | 0.021 | 0.025 |
| w/o $L_{\text{reg}}$ | 0.002 | 0.004 | 0.014 | 0.022 |
| Ours | 0.001 | 0.002 | 0.012 | 0.020 |

| Dataset | Internal | | | External | | |
|---|---|---|---|---|---|---|
| Metric | C | R | L | C | R | L |
| Notewise | 0.782 | 0.916 | 0.632 | 0.775 | 0.914 | 0.656 |
| CVAE | 0.798 | 0.812 | 0.620 | 0.773 | 0.802 | 0.649 |
| $L_{\text{pln}}$ | 0.693 | 0.852 | 0.323 | 0.694 | 0.834 | 0.324 |
| $L_{\text{pln}} + L_{\text{str}}$ | 0.633 | 0.882 | 0.253 | 0.639 | 0.865 | 0.277 |
| w/o $L_{\text{fac}}$ | 0.831 | 0.846 | 0.789 | 0.832 | 0.831 | 0.847 |
| w/o $L_{\text{reg}}$ | 0.804 | 0.955 | 0.653 | 0.808 | 0.946 | 0.657 |
| Ours | 0.942 | 0.953 | 0.976 | 0.944 | 0.945 | 0.977 |

- **C**: Consistency of the controlled attribute. / **R**: Restrictiveness of the uncontrolled attribute.
- **L**: Linearity b/t the controlled attribute and corresponding latent dimension.
- Our model shows the best scores in most metrics for disentanglement & controllability.
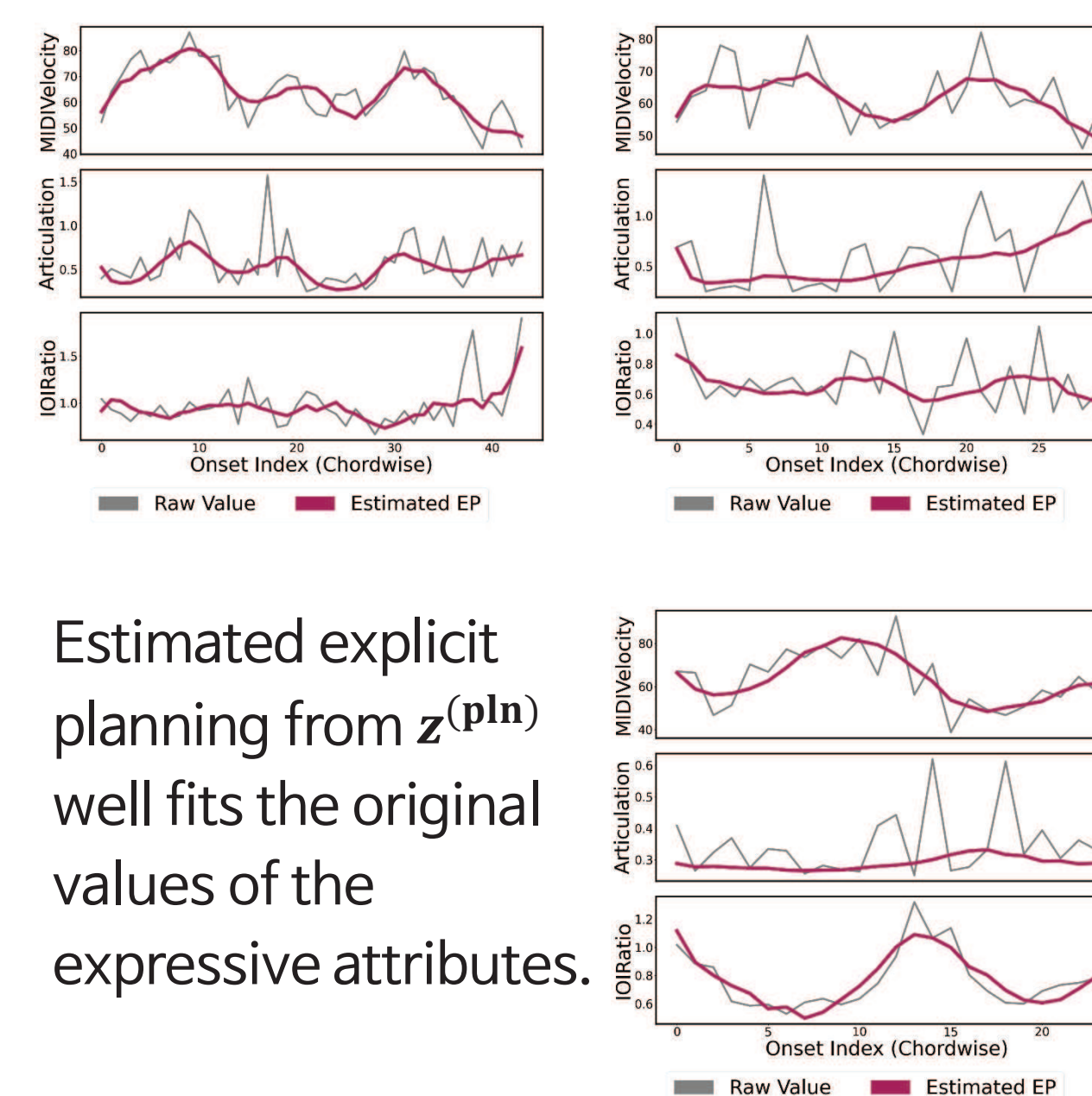
### III. Listening Test

| Metric | Winning Rate (Human-likeness) | | |
|---|---|---|---|
| Group | T | UT | Overall |
| Notewise | 0.317(±0.223) | 0.541(±0.316) | 0.493(±0.309) |
| CVAE | 0.467(±0.356) | 0.477(±0.342) | 0.475(±0.338) |
| Ours | 0.417(±0.256) | 0.555(±0.256) | 0.525(±0.258) |



- Winning rate: a rate of winning plain MIDI (A/B test)
- Top-ranking rate: a rate of being the highest rank in winning rate.
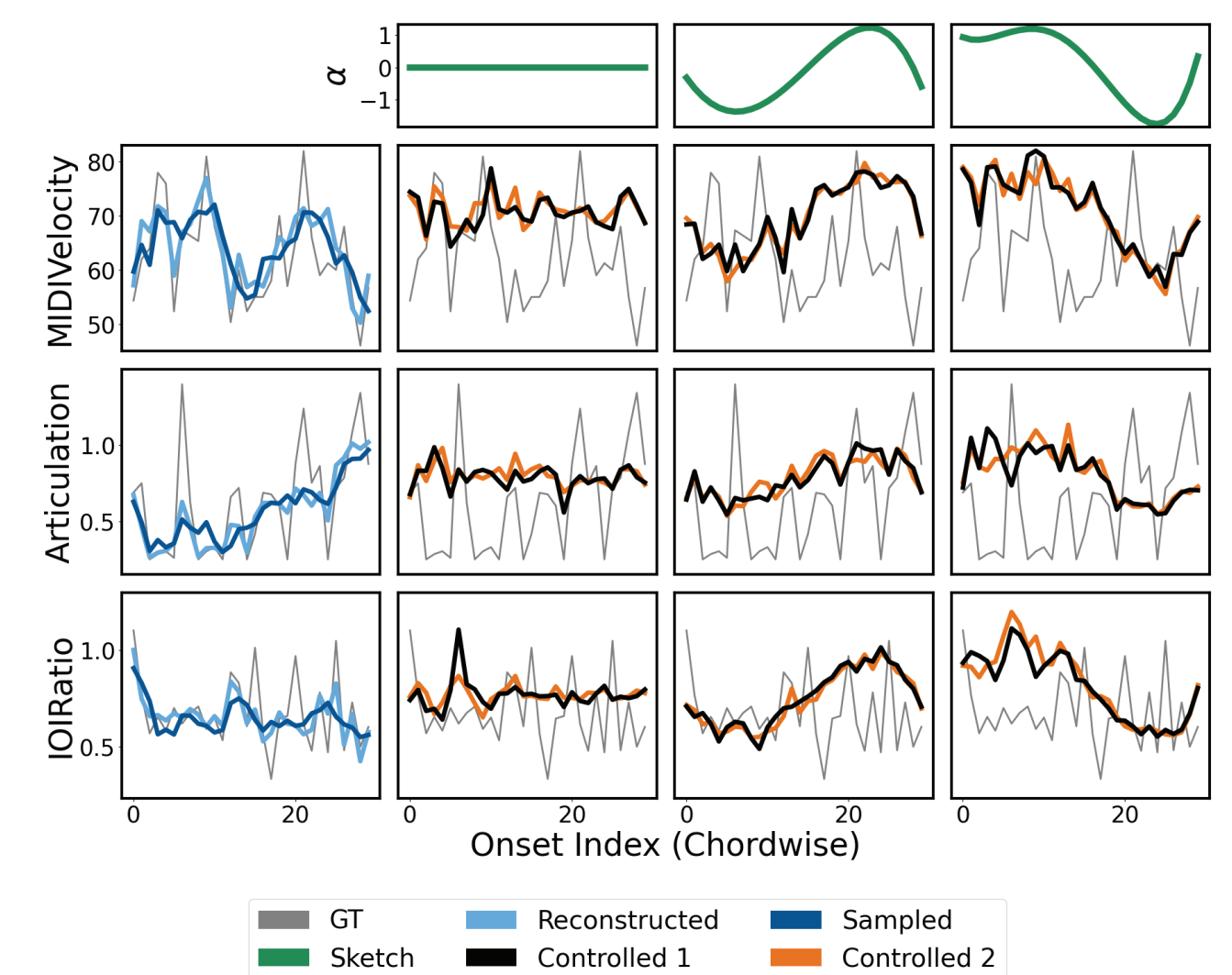- T/UT: musically trained (6) / untrained (22)

### III. Quantitative Results

**1) Estimating explicit planning**



- Estimated explicit planning from $z^{(\text{pln})}$ well fits the original values of the expressive attributes.

**2) Controlling expression with "sketches"**



GT | Reconstructed | Sampled
Sketch | Controlled 1 | Controlled 2

- $\alpha$: a sequence of values ("sketches") fed to a latent dimension $z^{(\text{pln})}$ for controlling a target expressive attribute.

## Conclusion

**Piano performance rendering with flexible musical expression**
- Our proposed system disentangles entire musical expression from piano performance and flexibly renders expressive piano performances in stable quality.
- Dynamics, articulation, and tempo can be independently controlled by our system while other structural attributes maintain their state.

**Future work**
- Deeper investigation for computing $I^{(\text{pln})}$ with other possible methods.
- Outputs can be rendered from scratch with random $z^{(\text{pln})}$ and $z^{(\text{str})}$. However, the random $z^{(\text{pln})}$ does not inherit temporal dependency without given sketch. Future study is needed for inferring $z^{(\text{pln})}$ that has temporal dependency without any specific sketch given as the input.