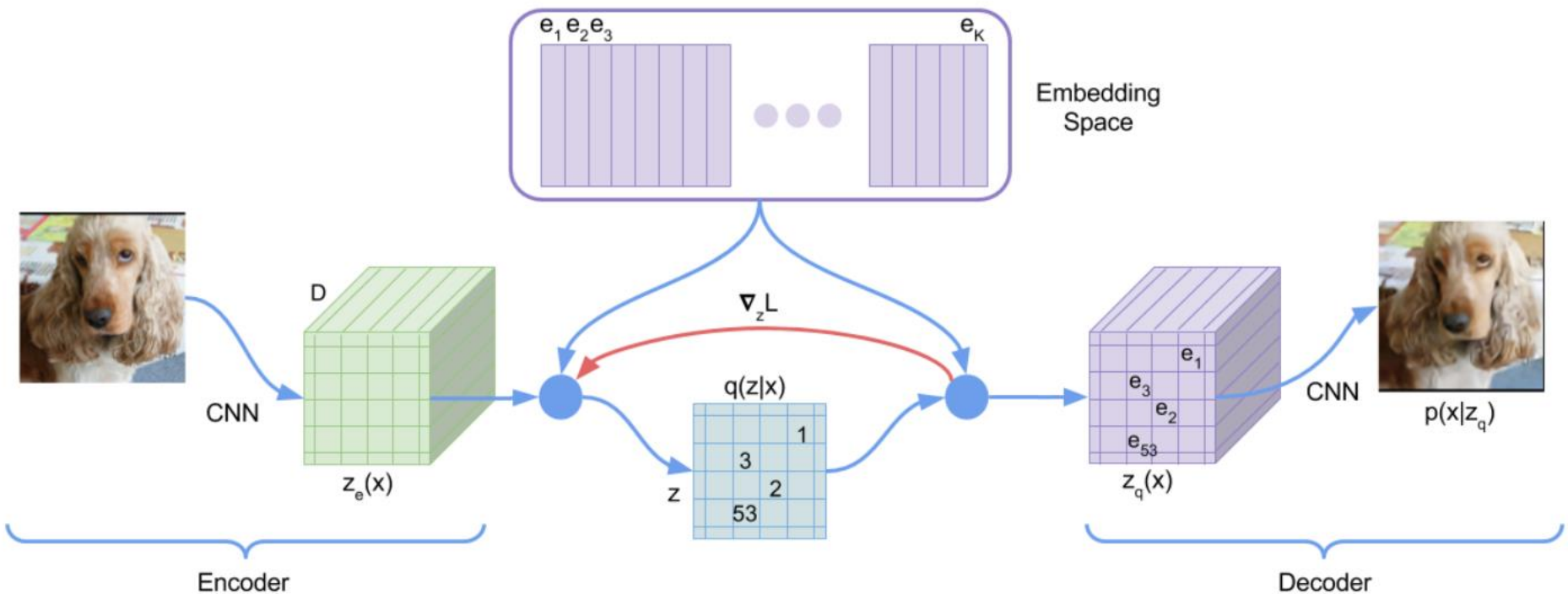# Jukebox: A Generative Model for Music

Speaker: Petr Grinberg

# Recall: VQ-VAE



New pictures can be generated from autoregressive distribution over z, p(z)

# Recall: VQ-VAE, Training Loss

$$\mathcal{L} = \mathcal{L}_{\text{recons}} + \mathcal{L}_{\text{codebook}} + \beta \mathcal{L}_{\text{commit}}$$
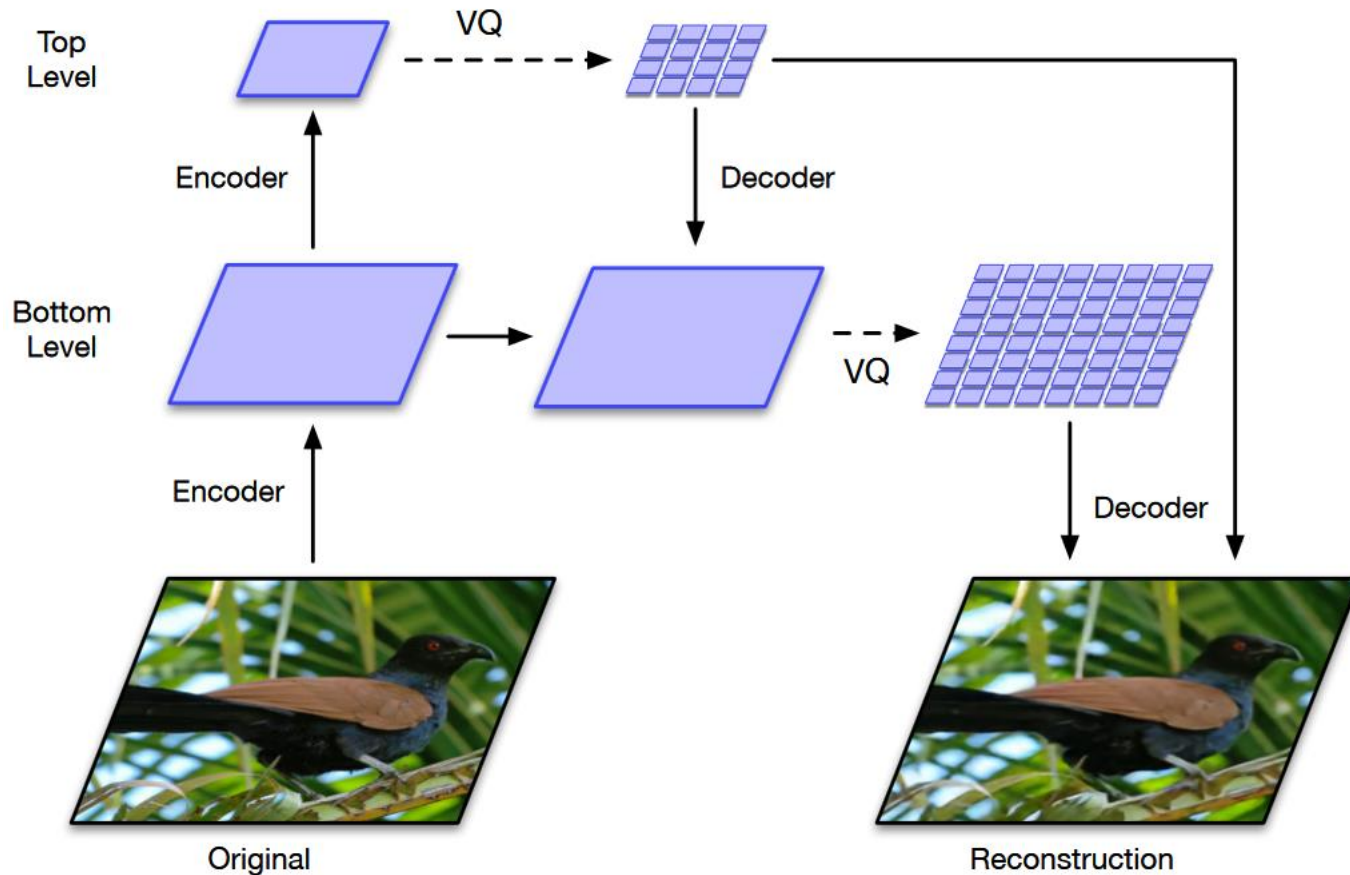
$$\mathcal{L}_{\text{recons}} = \frac{1}{T} \sum_t \|\mathbf{x}_t - D(\mathbf{e}_{z_t})\|_2^2$$

$$\mathcal{L}_{\text{codebook}} = \frac{1}{S} \sum_s \|\text{sg}[\mathbf{h}_s] - \mathbf{e}_{z_s}\|_2^2$$

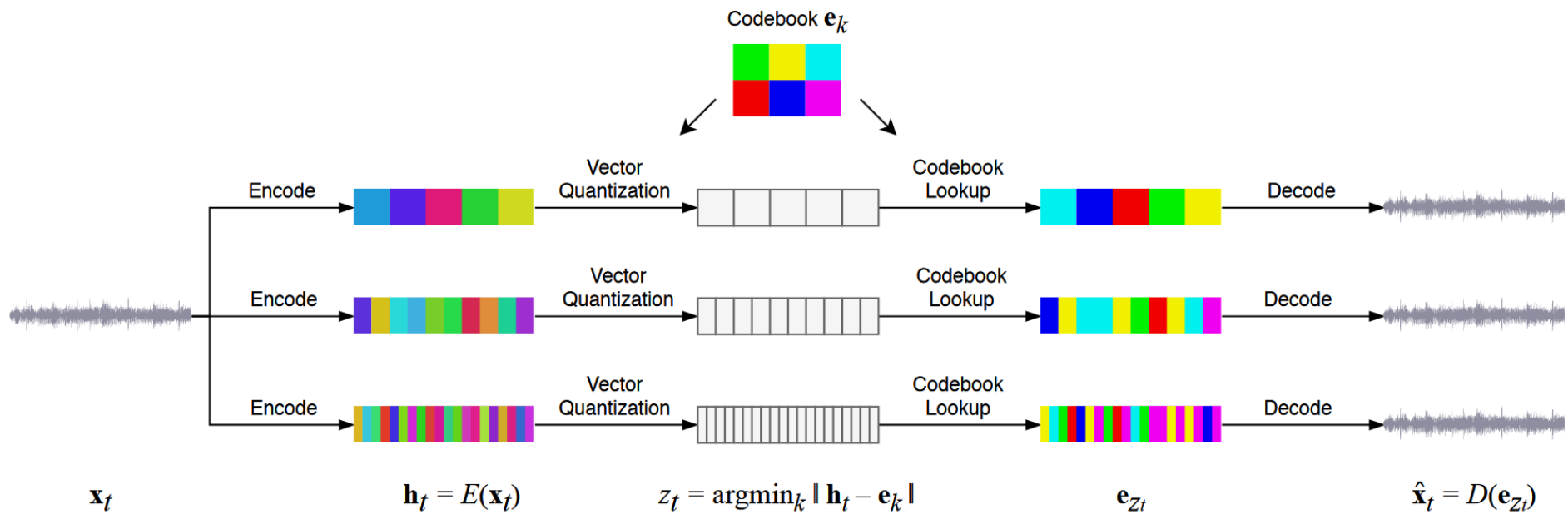$$\mathcal{L}_{\text{commit}} = \frac{1}{S} \sum_s \|\mathbf{h}_s - \text{sg}[\mathbf{e}_{z_s}]\|_2^2$$

Reconstruction Loss can vary, this one will be used in Jukebox
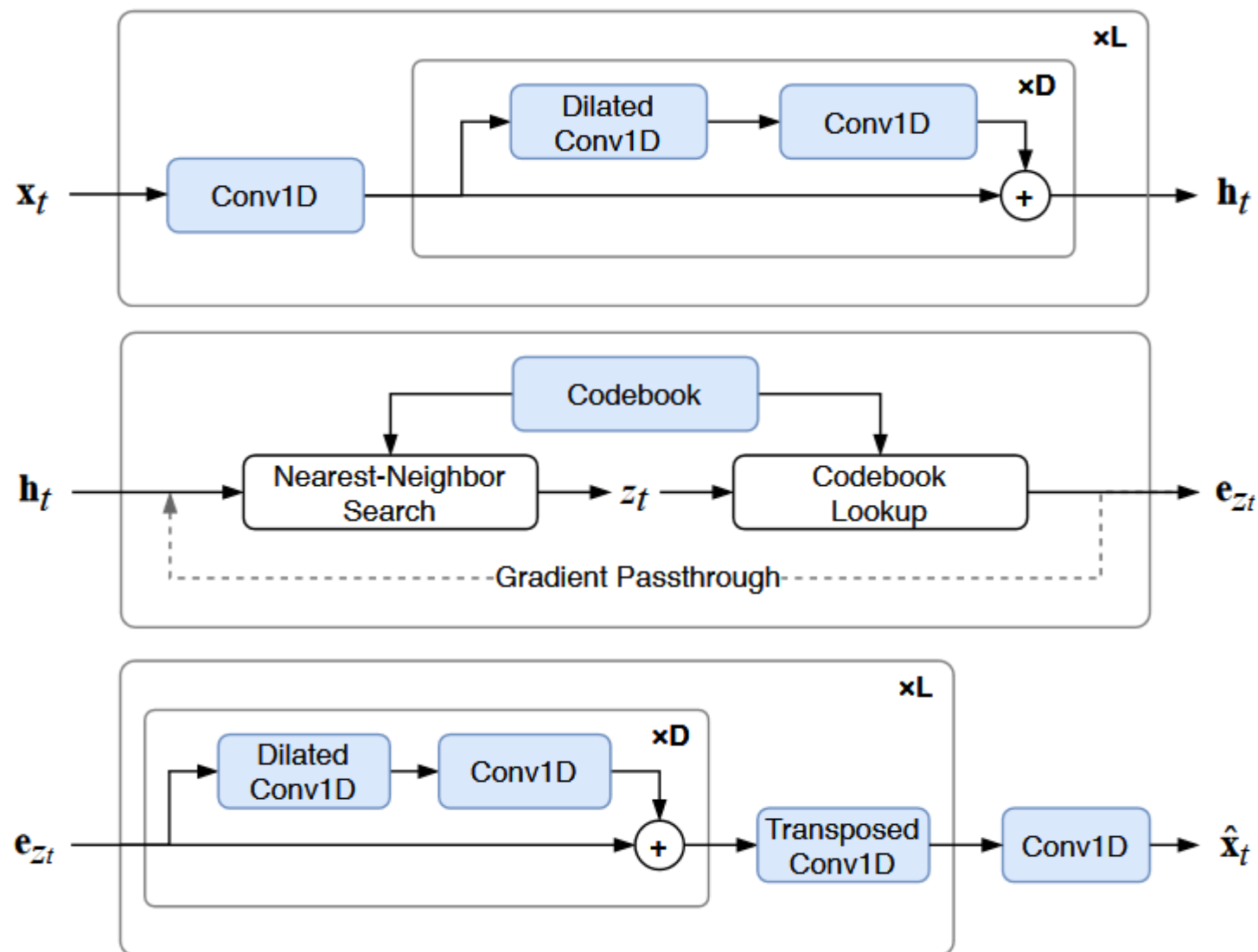
# Recall: VQ-VAE2 (Hierarchical)



Bottom Level: local features, Top Level: high-level semantics

# Jukebox: VQ-VAE
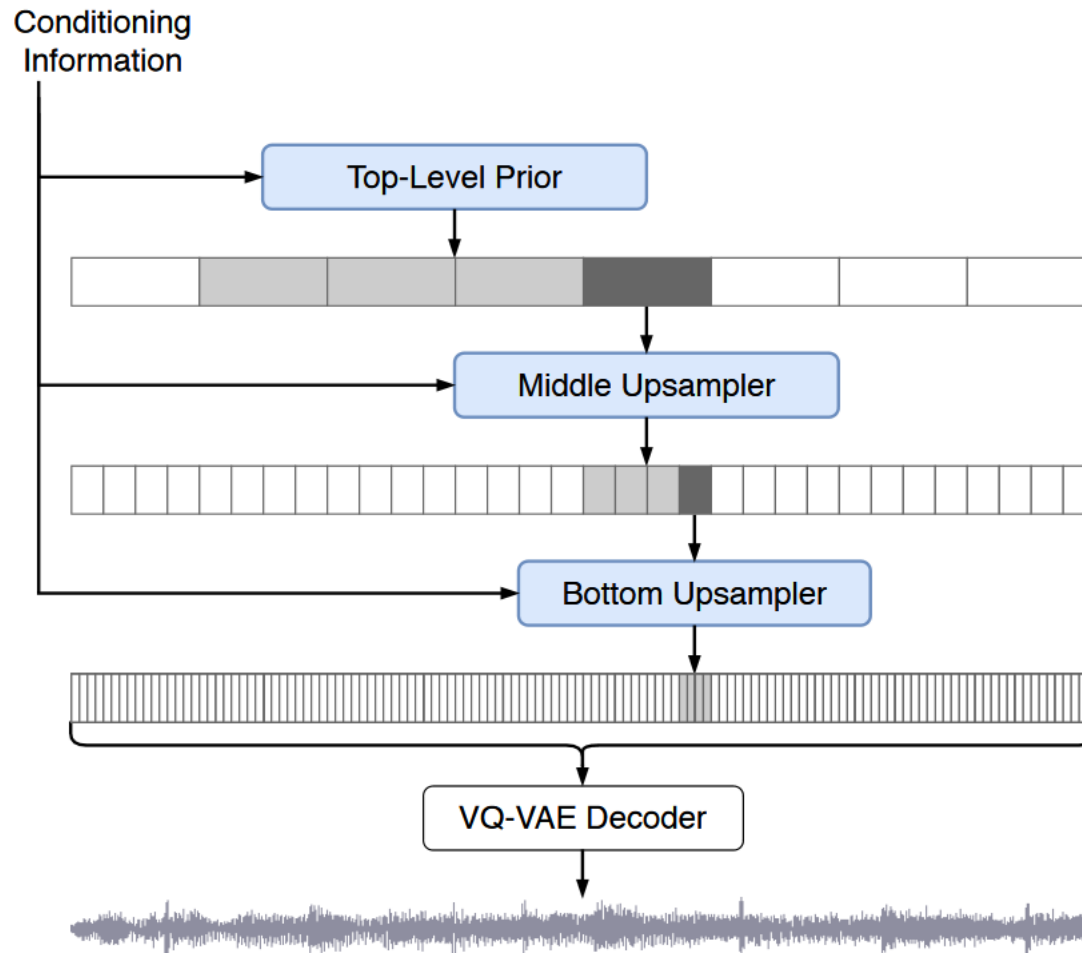
# Jukebox: VQ-VAE

# Jukebox: VQ-VAE, Training Hacks

- Random restarts for embeddings + EMA

- Separated Autoencoders

- Spectral Loss (for high-frequencies reconstruction):

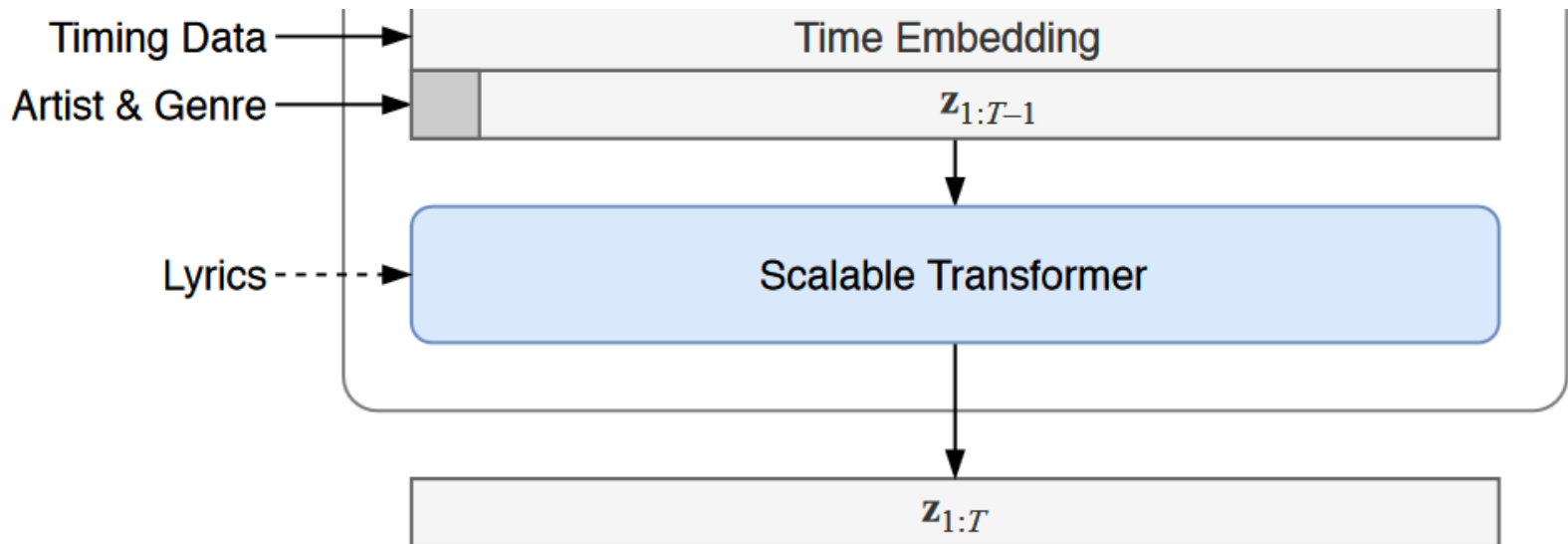$$\mathcal{L}_{\text{spec}} = \left\| |\text{STFT}(\mathbf{x})| - |\text{STFT}(\widehat{\mathbf{x}})| \right\|_2$$

# Jukebox: Learning Prior

$$p(\mathbf{z}) = p(\mathbf{z}^{\text{top}}, \mathbf{z}^{\text{middle}}, \mathbf{z}^{\text{bottom}})$$

$$= p(\mathbf{z}^{\text{top}}) p(\mathbf{z}^{\text{middle}} | \mathbf{z}^{\text{top}}) p(\mathbf{z}^{\text{bottom}} | \mathbf{z}^{\text{middle}}, \mathbf{z}^{\text{top}})$$
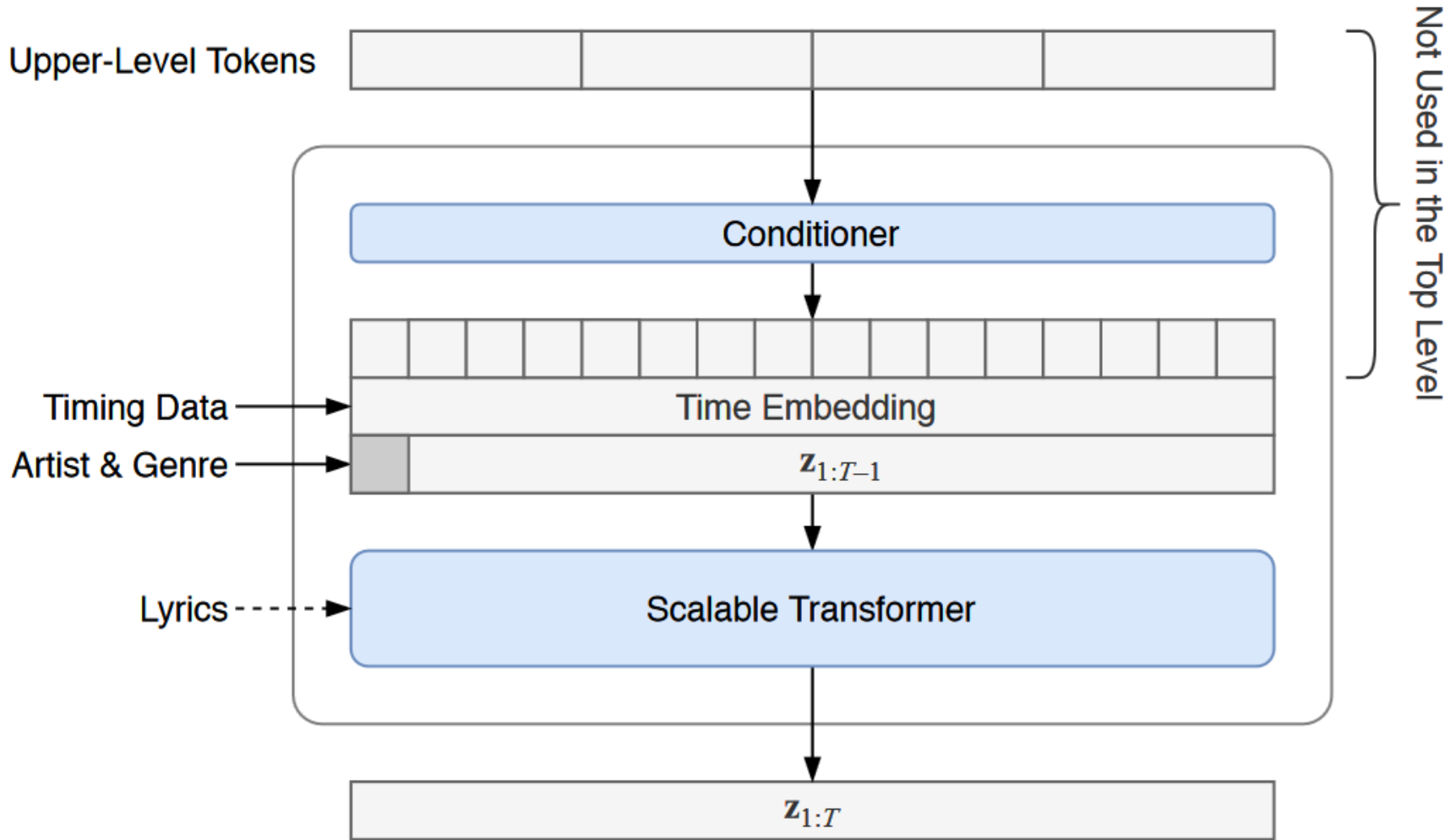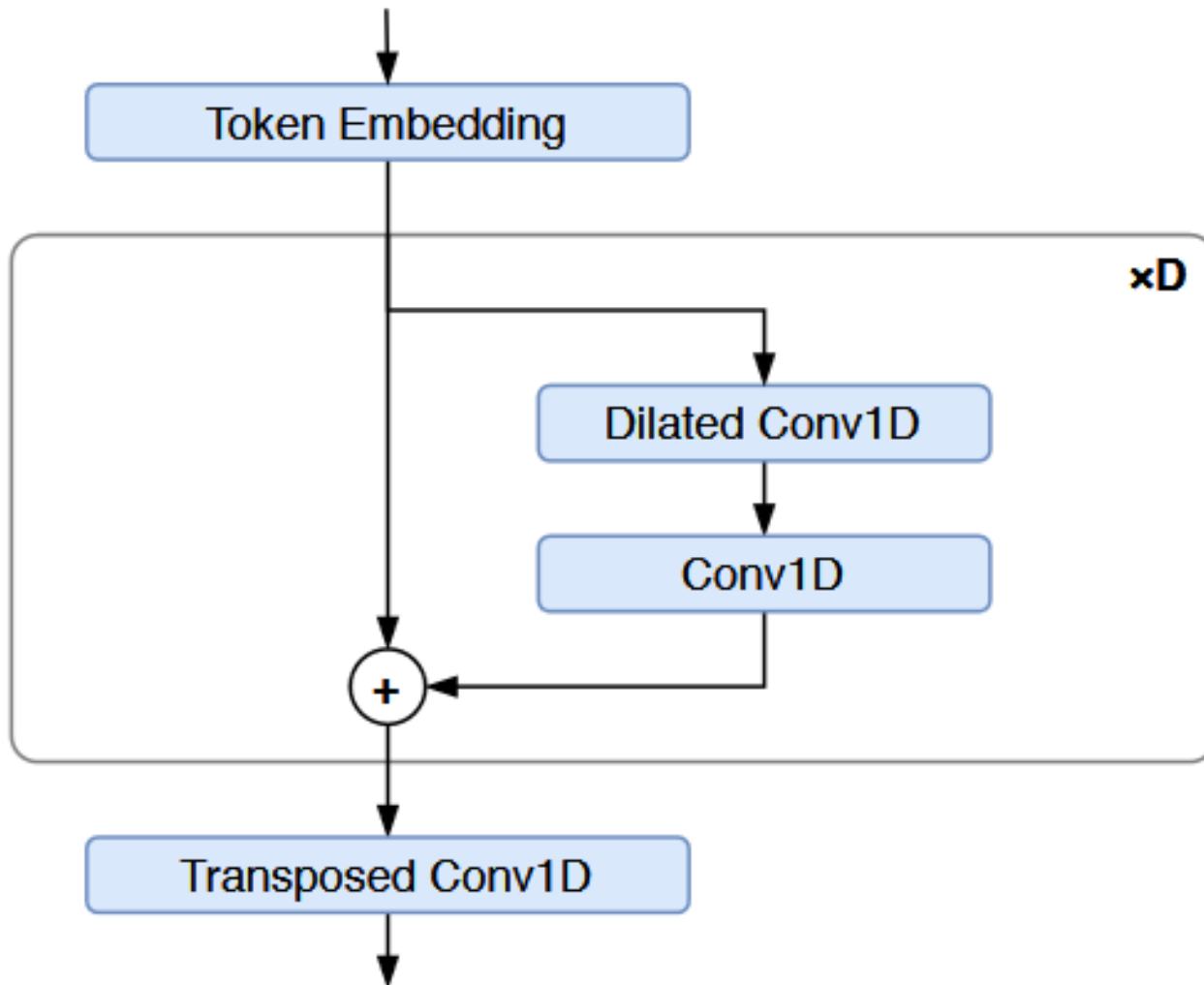
# Jukebox: Learning Prior

Top-level prior architecture:
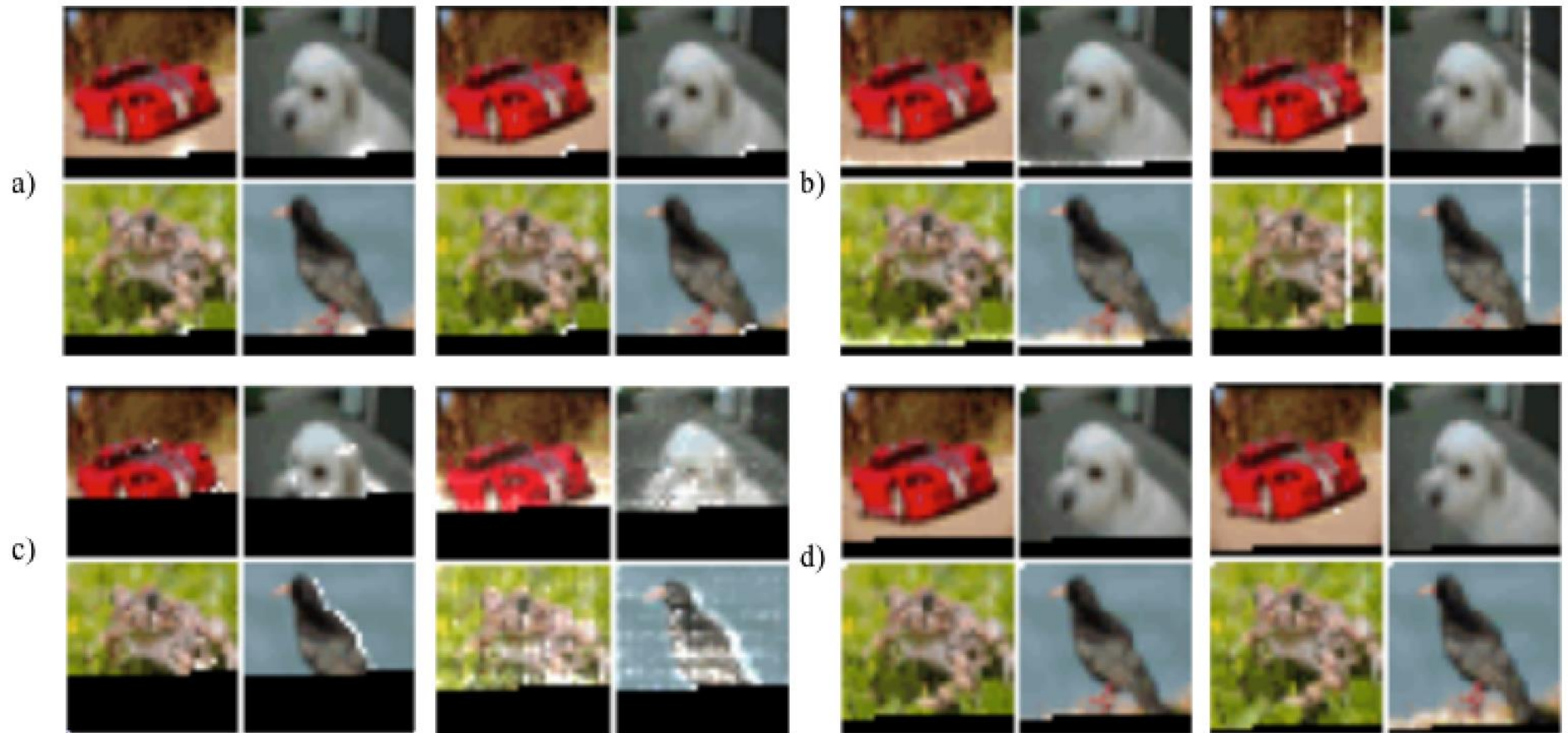
# Jukebox: Learning Prior
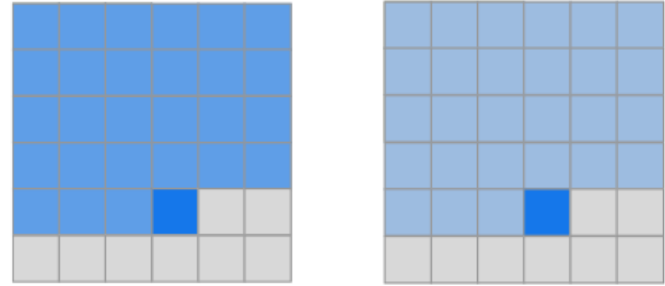
# Jukebox: Conditioner

# Jukebox: Scalable Transformer

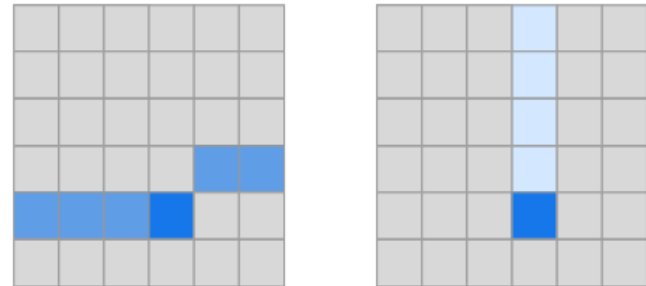Firstly, let's recall Sparse Transformer:

# Jukebox: Scalable Transformer
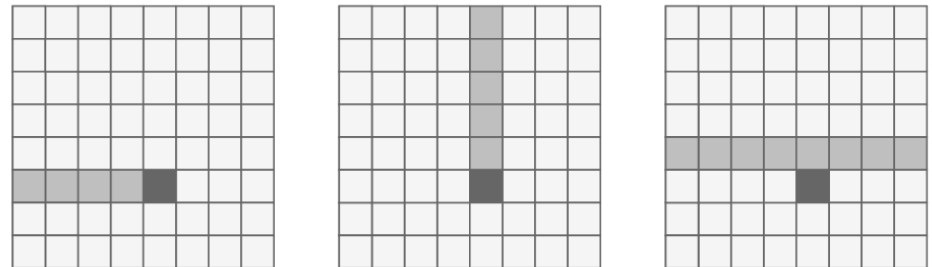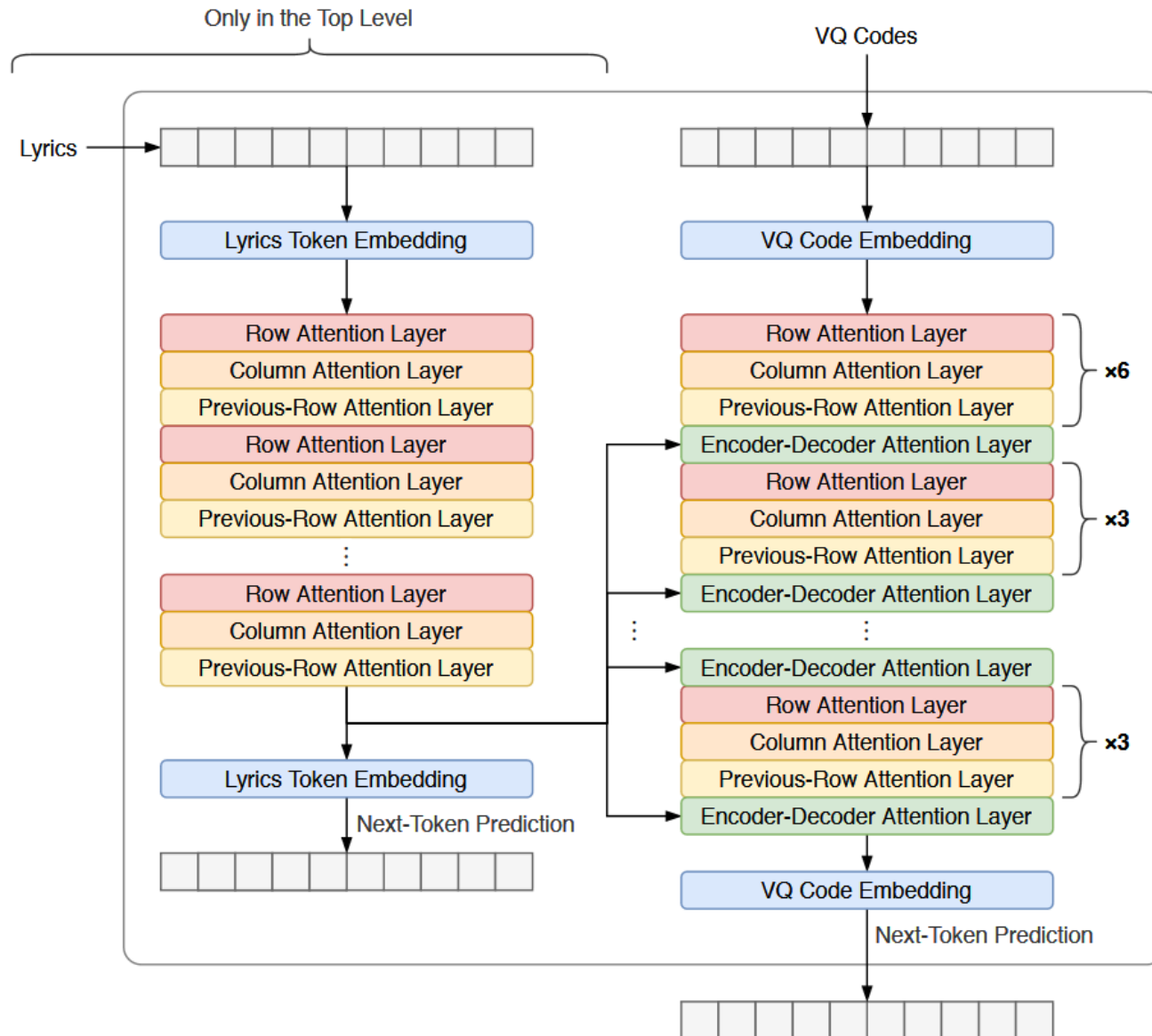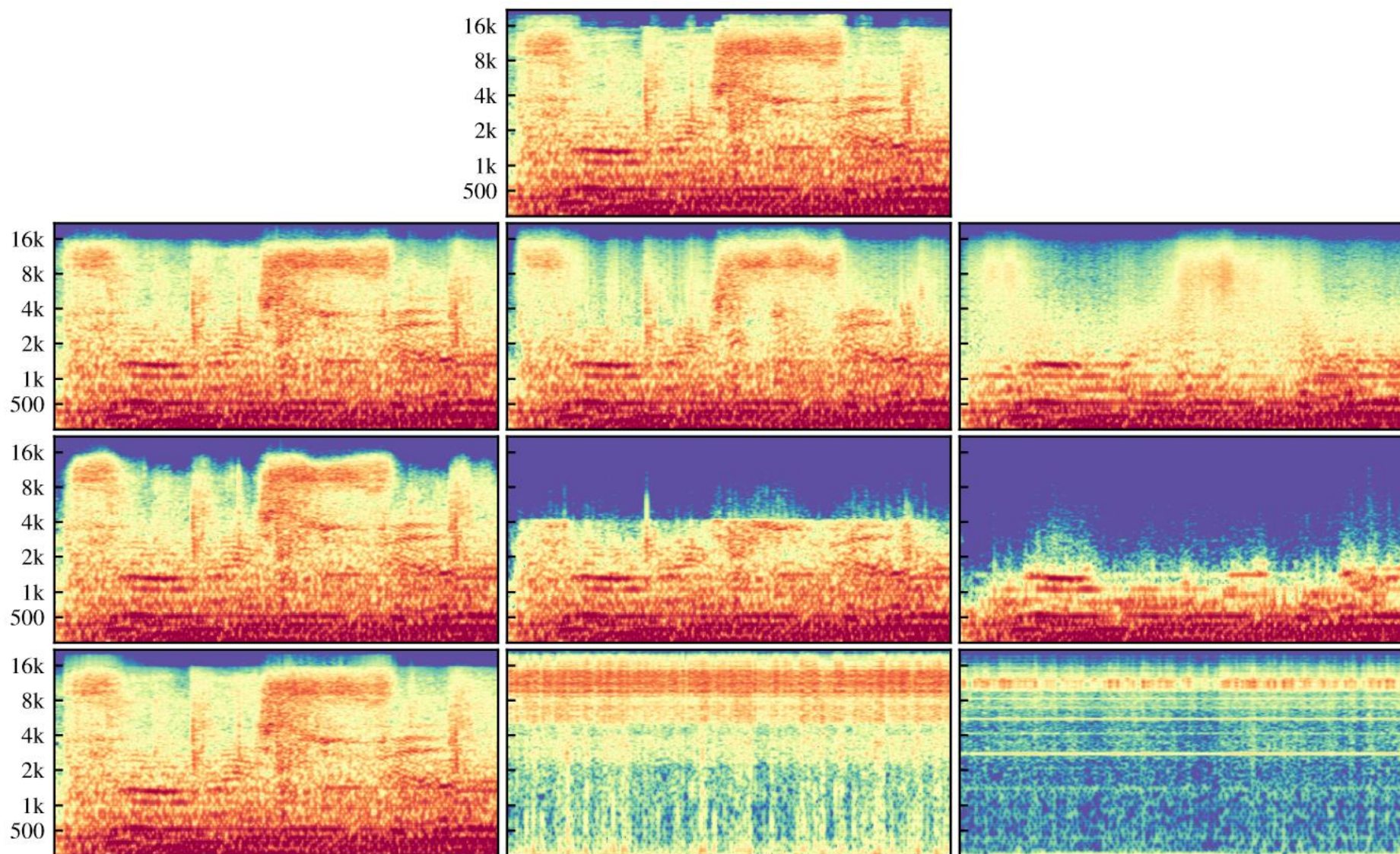
- Vanilla Transformer

- Sparse Transformer

- Scalable Transformer

# Jukebox: Lyrics Conditioning

# Jukebox: Ablation Study

# Jukebox: Ablation Study

| Level | Hop length | Spectral convergence (dB) Without restart | With restart |
|---|---|---|---|
| Bottom | 8 | −21.1 | −23.0 |
| Middle | 32 | −12.4 | −12.4 |
| Top | 128 | −8.3 | −8.3 |

| Codebook size | Spectral convergence (dB) |
|---|---|
| 256 | −15.9 |
| 2048 | −23.0 |
| No quantization | −40.5 |

| Ablation | Spectral convergence (dB) |
|---|---|
| None | −8.3 |
| Without spectral loss | −6.3 |
| With single autoencoder | 2.9 |