

WAVEGLOW: A FLOW-BASED GENERATIVE NETWORK FOR SPEECH SYNTHESIS

Ryan Prenger, Rafael Valle, Bryan Catanzaro

NVIDIA Corporation

Reviewer: Mark Baushenko



Glow + WaveNet = WaveGlow

INTRODUCTION

PROBLEM

WAVEGLOW

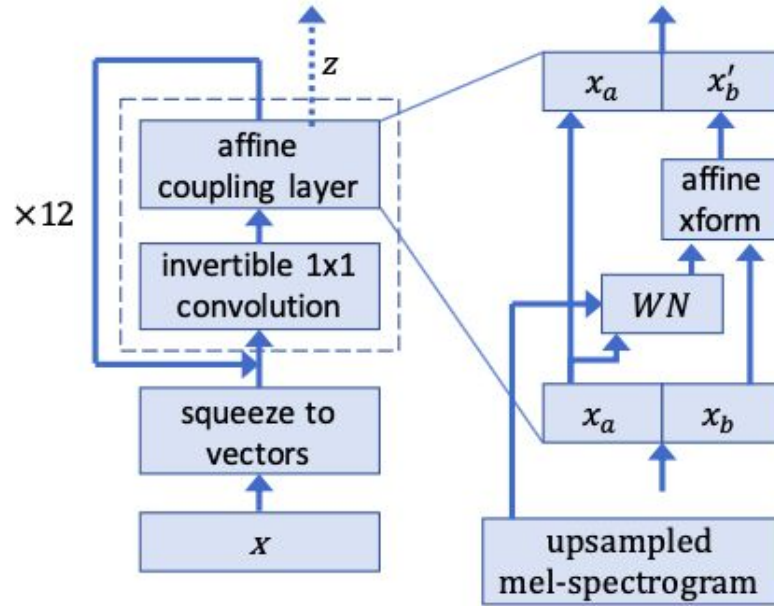


Fig. 1: WaveGlow network

1x1 Invertible Convolution

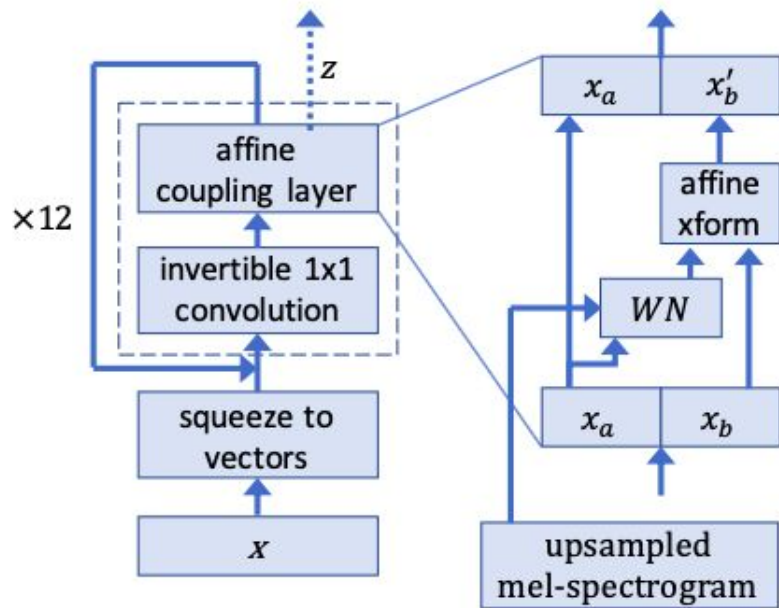


Fig. 1: WaveGlow network

$$f_{conv}^{-1} = \mathbf{W}x$$

$$\log |\det(\mathbf{J}(f_{conv}^{-1}(x)))| = \log |\det \mathbf{W}|$$

1x1 Invertible Convolution

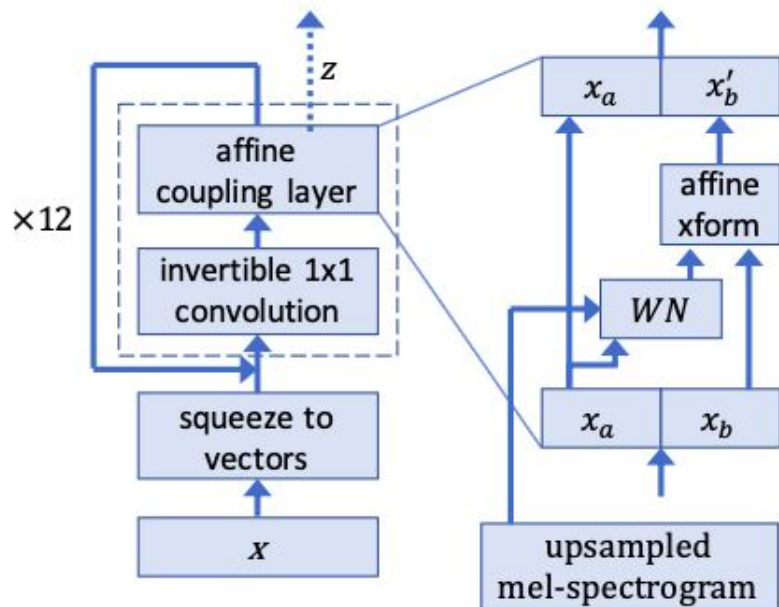


Fig. 1: WaveGlow network

$$f_{conv}^{-1} = \mathbf{W} \mathbf{x}$$

$$\log |\det(\mathbf{J}(f_{conv}^{-1}(\mathbf{x})))| = \log |\det \mathbf{W}|$$

$$\log p_{\theta}(\mathbf{x}) = -\frac{\mathbf{z}(\mathbf{x})^T \mathbf{z}(\mathbf{x})}{2\sigma^2} + \sum_{j=0}^{\#coupling} \log s_j(\mathbf{x}, mel-spectrogram) + \sum_{k=0}^{\#conv} \log \det |\mathbf{W}_k|$$

Audio quality comparison

Model	Mean Opinion Score (MOS)
Griffin-Lim	3.823 ± 0.1349
WaveNet	3.885 ± 0.1238
WaveGlow	3.961 ± 0.1343
Ground Truth	4.274 ± 0.1340

Thanks for attention