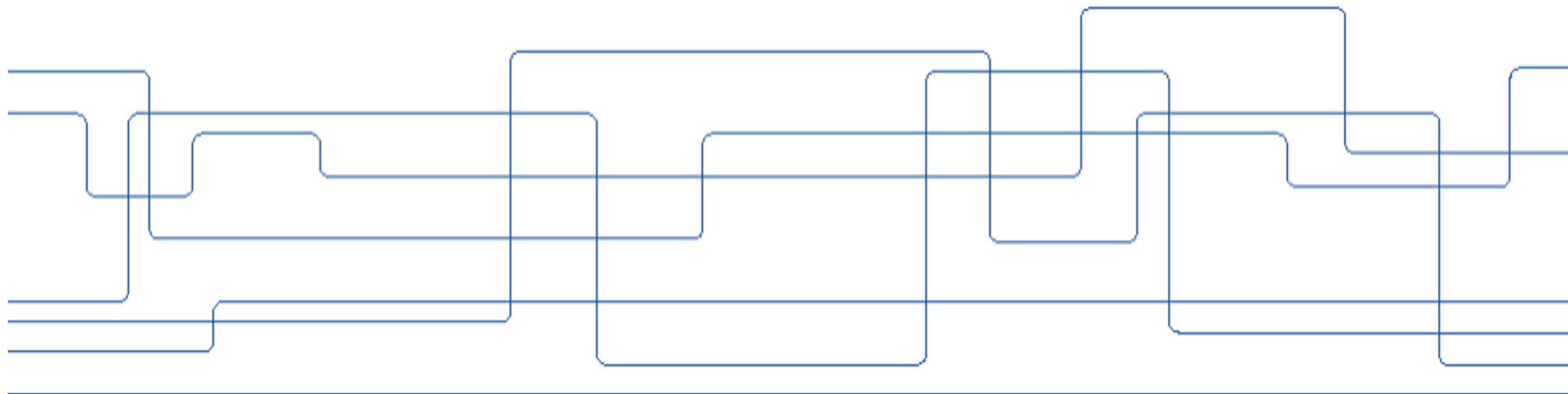


# Full-Glow: Fully conditional Glow for more realistic image generation

Moein Sorkhei, Gustav Eje Henter, Hedvig Kjellström





# Motivation

Synthesizing images of street scenes e.g. for training autonomous cars:

- Approach 1: Using graphic engines
  - Requires a lot of hand-engineering of features
  - Sometimes the synthesized images are not natural
- Approach 2: Using generative models
  - Ideally they could synthesize diverse images that resemble natural images
  - Image-to-image translation is commonly used to synthesize better target images

Popular generative models:

- GANs are extensively used thanks to their ability in synthesizing sharp images
- Flow-based models
  - Flexible probabilistic models with exact likelihood evaluation thanks to invertibility
  - After the publication of Glow, they have gained more attraction



# Contributions

Prior work using Glow for image-to-image translation:

- C-Glow: Makes the Glow layers conditional for image generation
- DUAL-Glow: Uses two Glows for medical image modality transfer, where the learned latent distribution of the target Glow is conditional
- C-Flow: Uses two Glows and makes all the coupling layers conditional

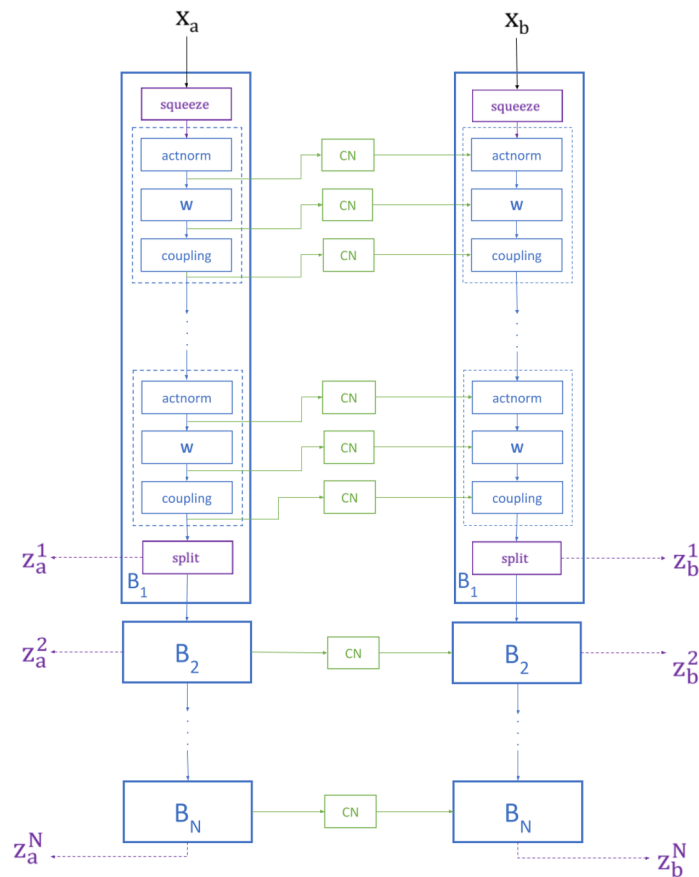
In our work, we build on the previous works by using two Glows and applying conditioning networks to all layers of the network, resulting in a fully conditional model.

- Lu, You, and Bert Huang. "Structured Output Learning with Conditional Generative Flows." *AAAI*. 2020.
- Sun, Haoliang, et al. "Dual-glow: Conditional flow-based generative model for modality transfer." *Proceedings of the IEEE International Conference on Computer Vision*. 2019.
- Pumarola, Albert, et al. "C-flow: Conditional generative flow models for images and 3d point clouds." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2020.

# Method

Full-Glow is a fully conditional architecture

$$\frac{1}{N} \left[ - \sum_{n=1}^N \lambda \log p_{\theta} \left( \mathbf{x}_a^{(n)} \right) - \sum_{n=1}^N \log p_{\phi} \left( \mathbf{x}_b^{(n)} \mid \mathbf{x}_a^{(n)} \right) \right]$$





# Setup of experiments

- Baseline models:
  - C-Glow: We consider two versions/configurations
    - Version 1: with deeper conditioning networks
    - Version 2: with deeper Glow (more flow blocks)
  - DUAL-Glow
  - Pix2pix: Conditional GAN for image-to-image translation
- We used the Cityscapes dataset
  - Natural street-scene images (RGB) with segmentation masks

- Lu, You, and Bert Huang. "Structured Output Learning with Conditional Generative Flows." *AAAI*. 2020.

- Sun, Haoliang, et al. "Dual-glow: Conditional flow-based generative model for modality transfer." *Proceedings of the IEEE International Conference on Computer Vision*. 2019.

- Isola, Phillip, et al. "Image-to-image translation with conditional adversarial networks." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017.

- Cordts, Marius, et al. "The cityscapes dataset for semantic urban scene understanding." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.



# Quantitative comparison

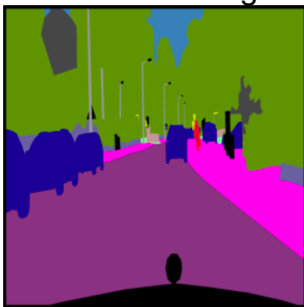
- Bits per dimension (BPD) for likelihood-based models
- For each model, we also resynthesized the validation set three times, conditioned on the segmentation masks
  - PSPNet segmentations of the generated images were then evaluated against the ground-truth segmentations in three ways

<b>Model</b>	<b>Cond. Mean BPD</b>	<b>Mean pixel acc.</b>	<b>Mean class acc.</b>	<b>Mean class IoU</b>
C-Glow v.1 [23]	2.568	$35.02 \pm 0.56$	$12.15 \pm 0.05$	$7.33 \pm 0.09$
C-Glow v.2 [23]	2.363	$52.33 \pm 0.46$	$17.37 \pm 0.21$	$12.31 \pm 0.24$
Dual-Glow [36]	2.585	$71.44 \pm 0.03$	$23.91 \pm 0.19$	$18.96 \pm 0.17$
pix2pix [15]	—	$60.56 \pm 0.11$	$22.64 \pm 0.21$	$16.42 \pm 0.06$
<b>Our model</b>	<b>2.345</b>	<b><math>73.50 \pm 0.13</math></b>	<b><math>29.13 \pm 0.39</math></b>	<b><math>23.86 \pm 0.30</math></b>
<i>Ground-truth</i>	—	<i>95.97</i>	<i>84.31</i>	<i>77.30</i>

Zhao, Hengshuang, et al. "Pyramid scene parsing network." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017.

# Visual comparison of model outputs

Conditioning



Ground truth



C-Glow v.1



C-Glow v.2



DUAL-Glow



pix2pix



Full-Glow sample 1

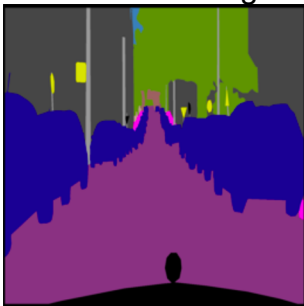


Full-Glow sample 2



# Visual comparison of model outputs

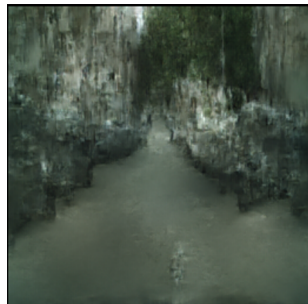
Conditioning



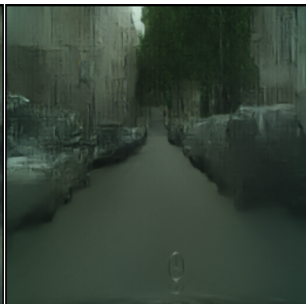
Ground truth



C-Glow v.1



C-Glow v.2



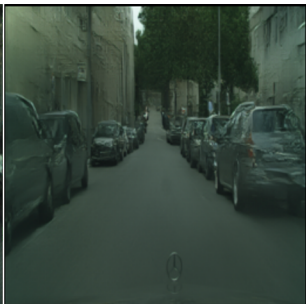
DUAL-Glow



pix2pix



Full-Glow sample 1

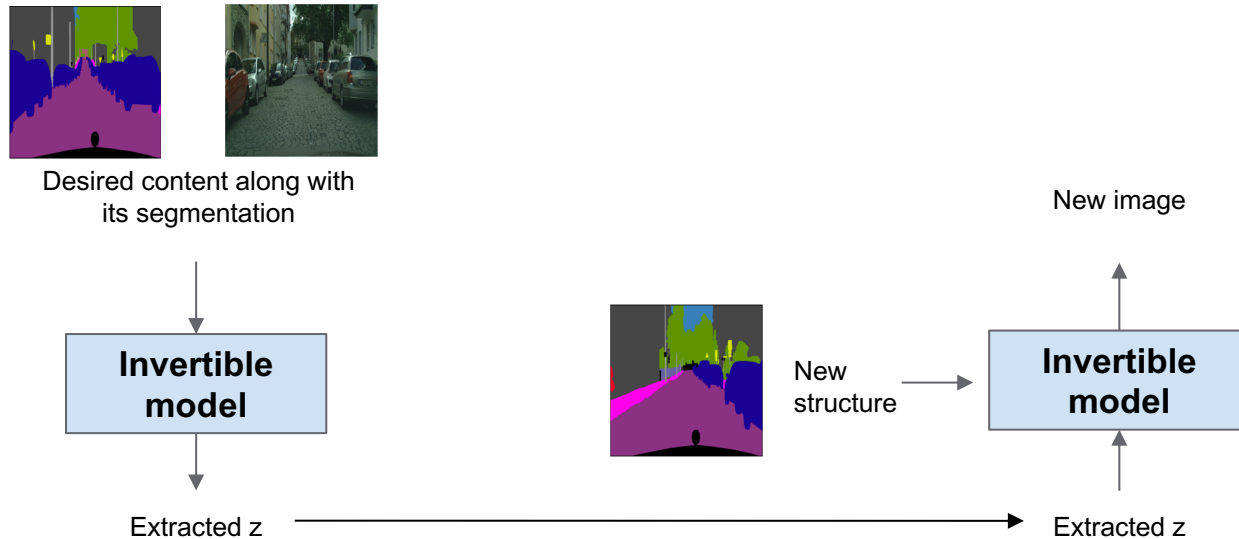


Full-Glow sample 2





# Using the model for content transfer

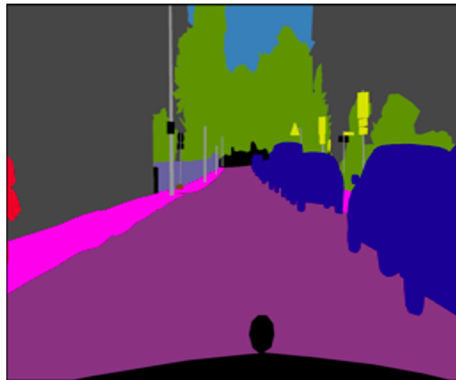


# Content transfer example

*Desired content*



*Desired structure*



*Content applied to structure*



Notice the difference in walls, asphalt, and the cars compared to ground-truth image

*Ground-truth image for the given structure* →



# Content transfer example

*Desired content*



*Desired structure*



*Content applied to structure*



Notice the difference in walls, asphalt, and the cars compared to ground-truth image

*Ground-truth image for the given structure* →





# Summary

- We proposed a fully conditional Glow-based architecture for more realistic conditional street-scene image generation
- Our improved conditioning allows for generating images that are more interpretable by a semantic classifier
- We synthesized higher-resolution images than previous works (see paper)
- We demonstrated promising results in content transfer
- Flow-based models can provide an alternative to GANs



Thank you for listening!