

Towards Object Shape Translation Through Unsupervised Generative Deep Models

Lies Bollens, Tinne Tuytelaars & José Oramas M.

KU Leuven, ESAT-PSI

September 25th, 2019



In a nutshell ...

What?

Translating shape, preserving style (colour, texture, etc)



In a nutshell ...

What?

Translating shape, preserving style (colour, texture, etc)



How?

- Learn shapes independently
- Learn how to map the learned shapes

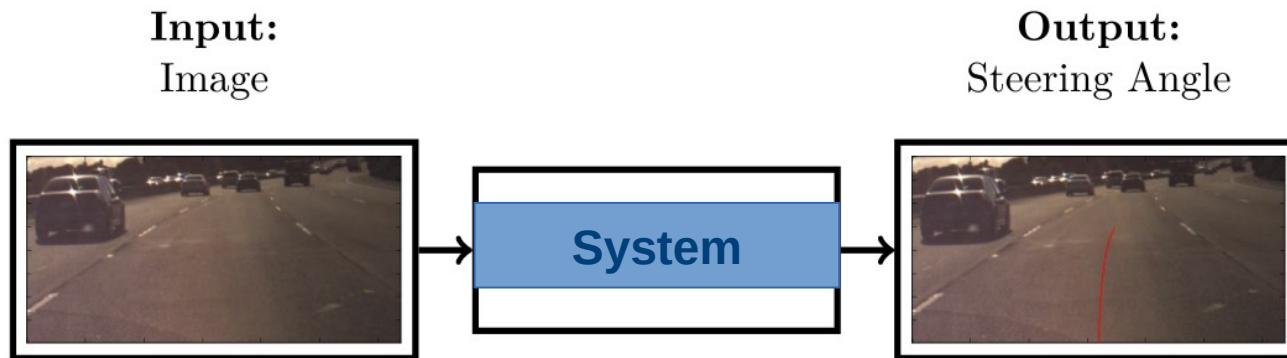
Some history

End-to-End Steering Prediction (2016 - 2017)

Previous Work ...

End-to-end Steering Prediction¹

Given an image sequence → Predict a steering angle



¹Heylen et al., "From Pixels to Actions: Learning to Drive a Car with Deep Neural Networks". WACV'18.

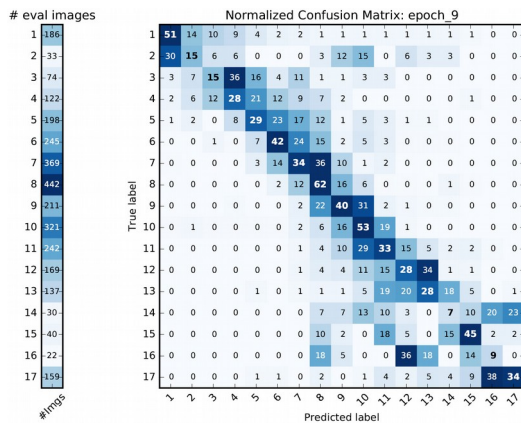
Previous Work ...

End-to-end Steering Prediction¹

Some Results

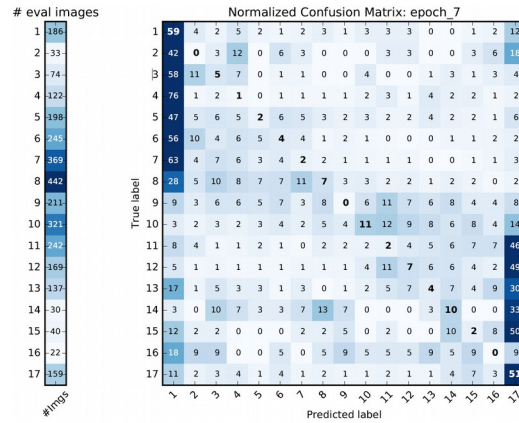
Data: Udacity & GTA-V Simulator

Regular only



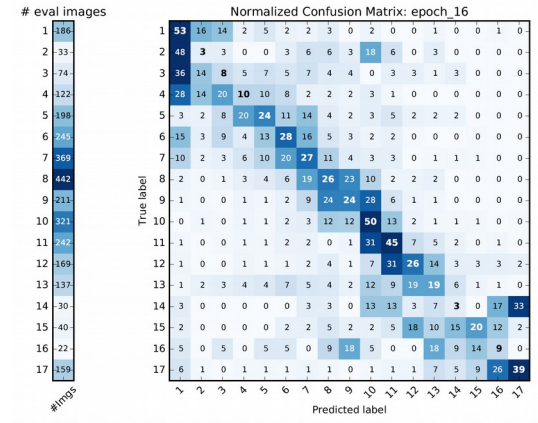
worst

Recovery only



best

Regular + Recovery



¹Heylen et al., "From Pixels to Actions: Learning to Drive a Car with Deep Neural Networks". WACV'18.

Nice, but ...

Training Conditions (simulator)



Testing Conditions (KITTI dataset)



Far from the testing conditions

Making simulation data more realistic

(2017 - 2018)

From the virtual world to reality

What?

Given simulated data → Bring it closer to testing conditions



Grand Theft Auto V



Example from the KITTI dataset

How?

- Domain/image translation
- Generative Adversarial Networks (GANs)¹

From GTA-V to realistic images



- Original image examples from GTA-V (*input*)
- GAN + Wasserstein Loss

Nice but...

Currently it is mostly about changing pixel colours,

What about the structure?

(shape of trees, models of cars, building architecture)

input



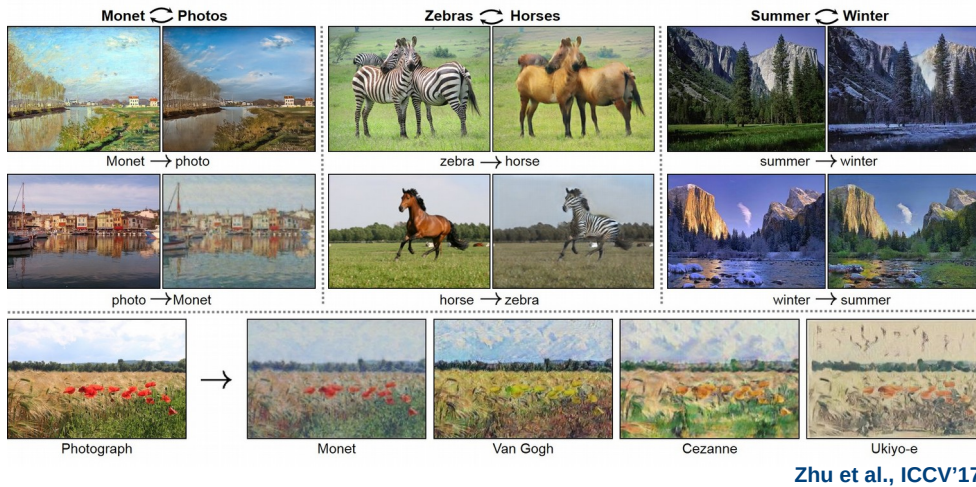
output



Translating Shapes, Preserving Style

Related Work

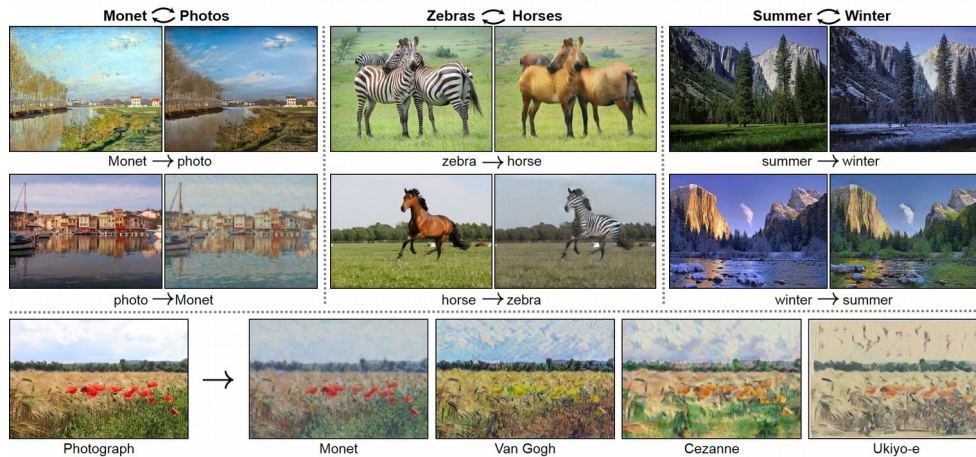
Neural Style Transfer



- Gatys et al., CVPR'16.
- Zhu et al., ICCV'17.
- Luan et al., CVPR'17.
- Mechrez et al., BMVC'17.
- Wang et al., CVPR'17
- Liu et al., NPAR'17,
- Dumoulin et al., ICLR'17.
- Chen et al., CVPR'17.
- Li et al. CVPR'17.
- Zhang et al., arXiv:1703.06953.
- Jing et al., ECCV'18.
- **Jing et al., arXiv:1705.04058.**

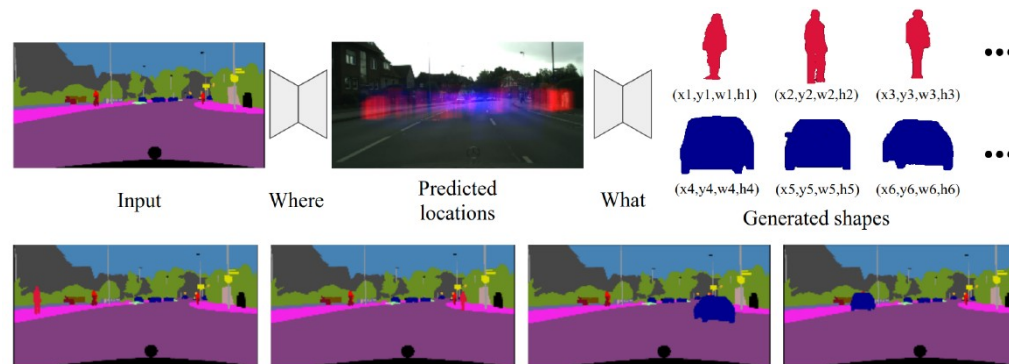
Related Work

Neural Style Transfer



Zhu et al., ICCV'17.

Shape Transfer



Lin et al., CVPR'18.

- Gatys et al., CVPR'16.
- Zhu et al., ICCV'17.
- Luan et al., CVPR'17.
- Mechrez et al., BMVC'17.
- Wang et al., CVPR'17
- Liu et al., NPAR'17,
- Dumoulin et al., ICLR'17.
- Chen et al., CVPR'17.
- Li et al. CVPR'17.
- Zhang et al., arXiv:1703.06953.
- Jing et al., ECCV'18.
- **Jing et al., arXiv:1705.04058.**

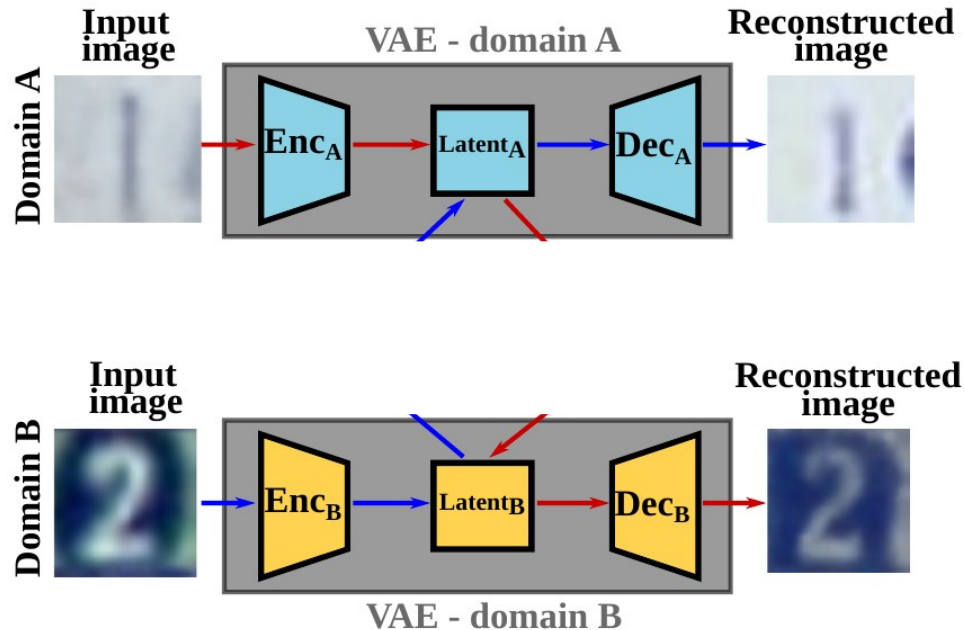
- Jaderberg et al., NIPS'15.
- Lee et al., NIPS'18.
- Lin et al., CVPR'18.

Proposed Method

Proposed Method

How? (notion)

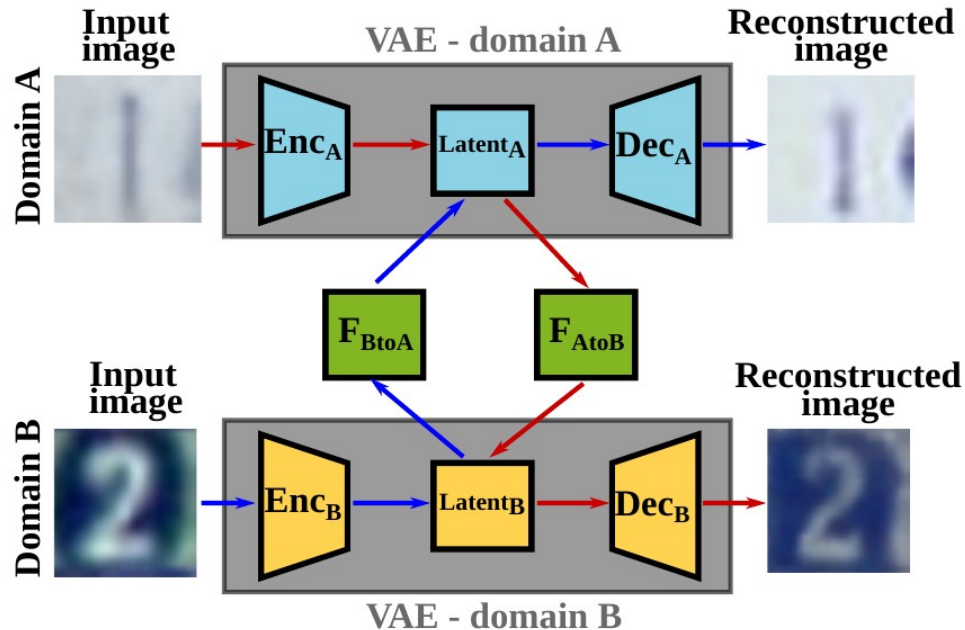
1- Learn how to model each domain independently.



Proposed Method

How? (notion)

- 1- Learn how to model each domain independently.
- 2- Learn a translation function between the domains¹.

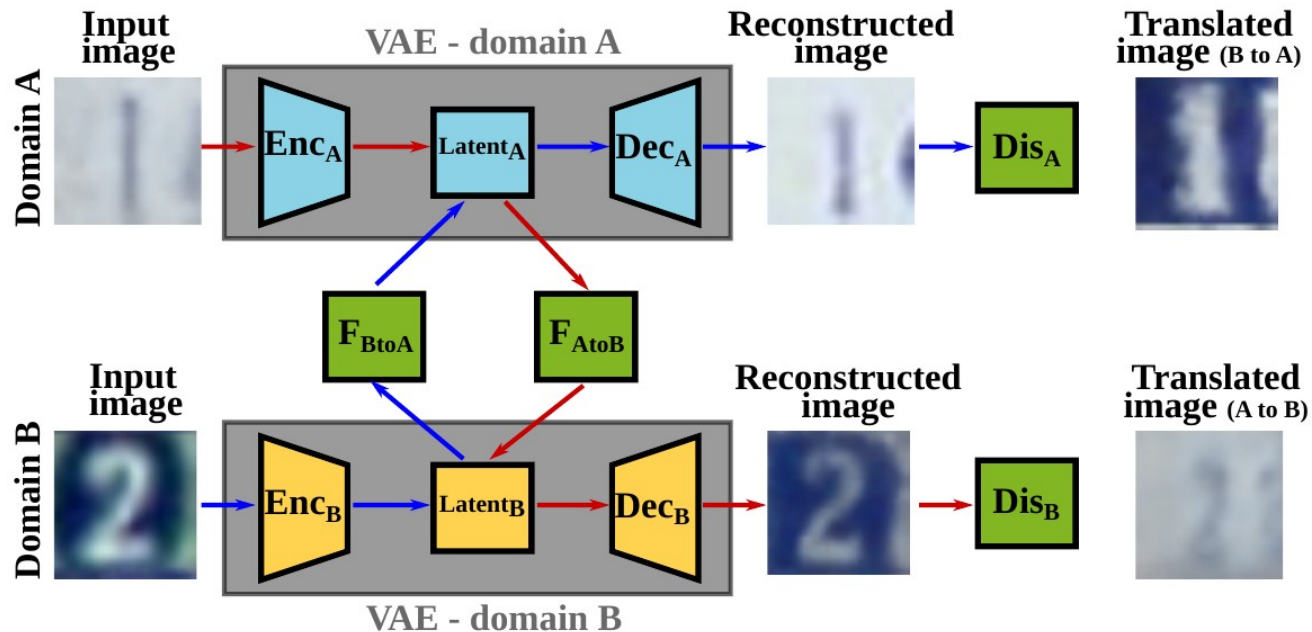


¹Lejeune et al. , "A Data Driven Similarity Measure and Example Mapping Function for General, Unlabelled Data Sets", ECAI'16.

Proposed Method

How? (notion)

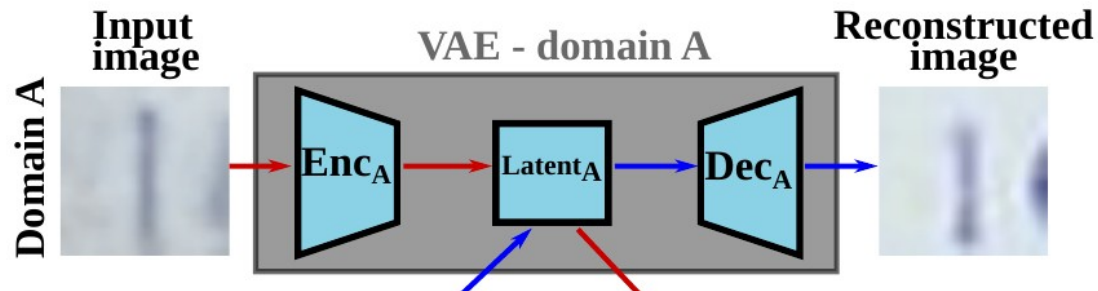
- 1- Learn how to model each domain independently.
- 2- Learn a translation function between the domains.



Proposed Method

Modelling Domain Information

- Learn how to model each domain independently.
 - Through variational autoencoders (VAE)¹

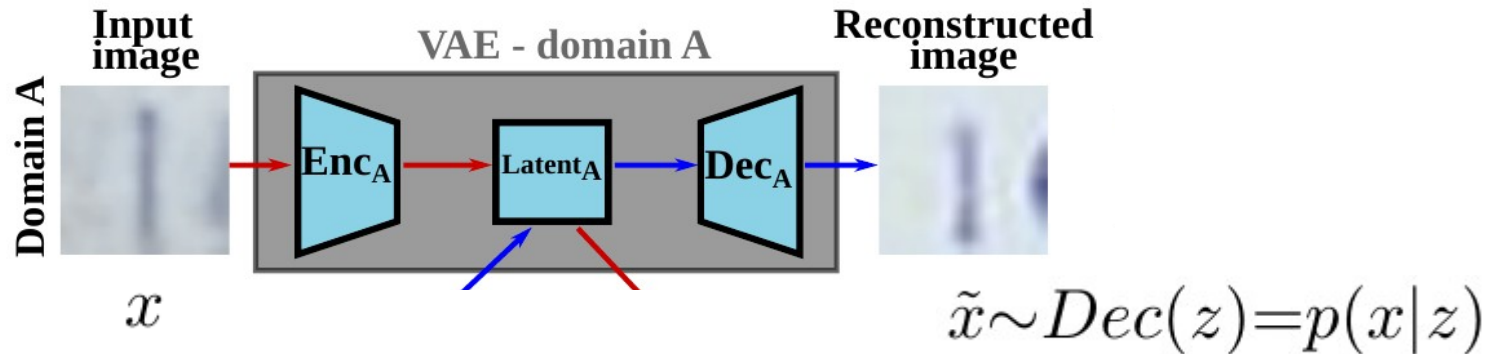


¹D. P Kingma and M. Welling, "Auto-Encoding Variational Bayes," in arXiv:1312.6114, 2013.

Proposed Method

Modelling Domain Information

- Learn how to model each domain independently.
 - Through variational autoencoders (VAE)¹

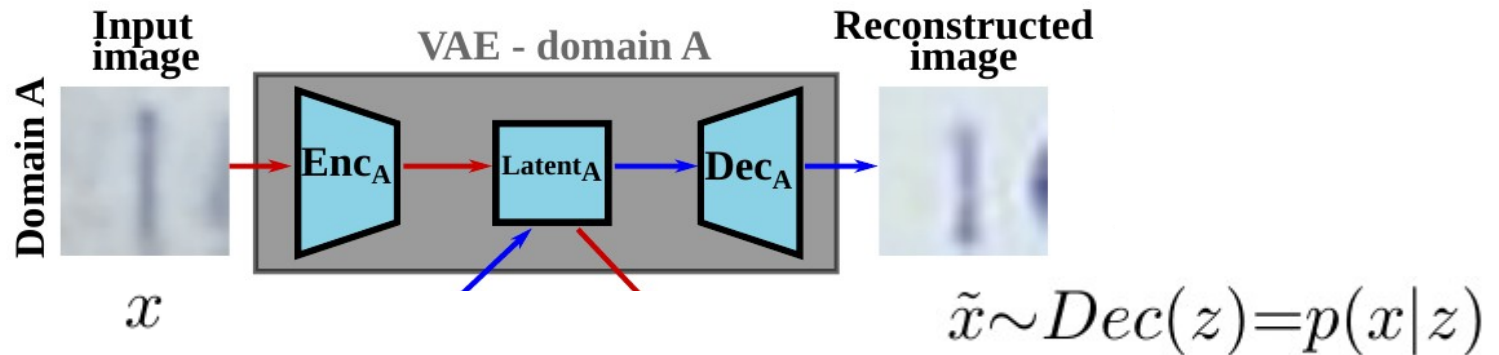


¹D. P Kingma and M. Welling, “Auto-Encoding Variational Bayes,” in arXiv:1312.6114, 2013.

Proposed Method

Modelling Domain Information

- Learn how to model each domain independently.
 - Through variational autoencoders (VAE)¹



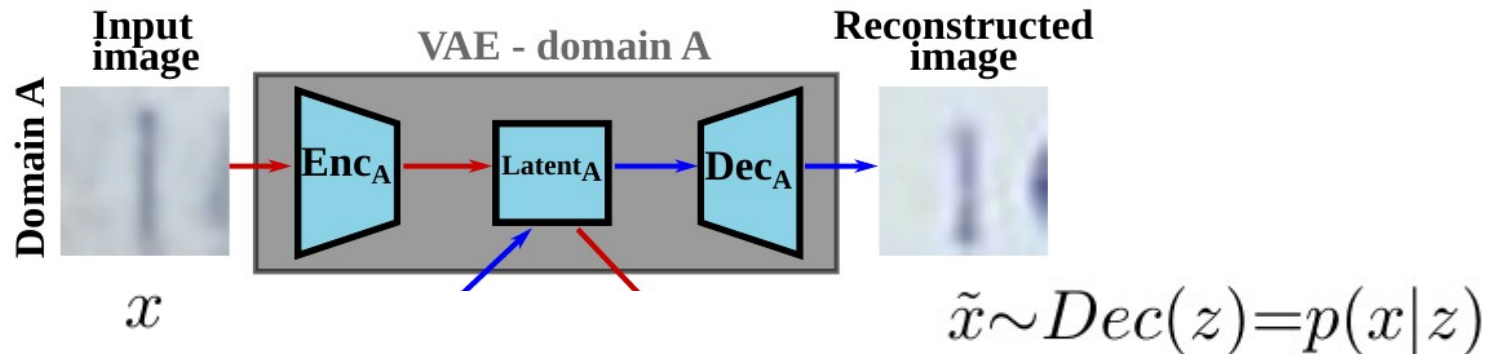
With latent space $\rightarrow z \sim Enc(x) = q(z|x)$

¹D. P Kingma and M. Welling, "Auto-Encoding Variational Bayes," in arXiv:1312.6114, 2013.

Proposed Method

Modelling Domain Information

- Learn how to model each domain independently.
 - Through variational autoencoders (VAE)¹



With latent space $\rightarrow z \sim Enc(x) = q(z|x)$

Applying the Loss

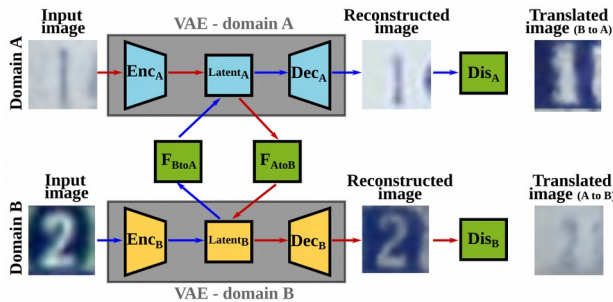
$$\mathcal{L}_{VAE} = - \mathbb{E}_{q(z|x)} [\log p(x|z)] + \mathcal{D}_{KL}(q(z|x) || p(z))$$

¹D. P Kingma and M. Welling, "Auto-Encoding Variational Bayes," in arXiv:1312.6114, 2013.

Proposed Method

Shape Translation

- Learn a mapping function between the [shape] domains.
 - Through an extended CycleGAN¹.



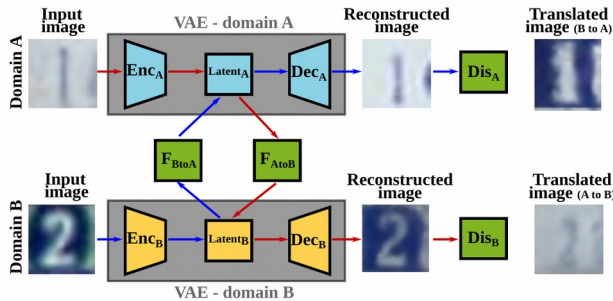
→ Applying the Loss

$$\begin{aligned}\mathcal{L} = & \mathcal{L}_{GAN}(Gen_A, Dis_A) + \mathcal{L}_{GAN}(Gen_B, Dis_B) \\ & + \lambda_{cycle} * \mathcal{L}_{cyc}(Gen_A, Gen_B) \\ & + \lambda_{sim} * \mathcal{L}_{sim}(x, y, Gen_A, Gen_B)\end{aligned}$$

Proposed Method

Shape Translation

- Learn a mapping function between the [shape] domains.
 - Through an extended CycleGAN¹.



→ Applying the Loss

$$\mathcal{L} = \mathcal{L}_{GAN}(Gen_A, Dis_A) + \mathcal{L}_{GAN}(Gen_B, Dis_B) \\ + \lambda_{cycle} * \mathcal{L}_{cyc}(Gen_A, Gen_B) \\ + \lambda_{sim} * \mathcal{L}_{sim}(x, y, Gen_A, Gen_B)$$

Where:

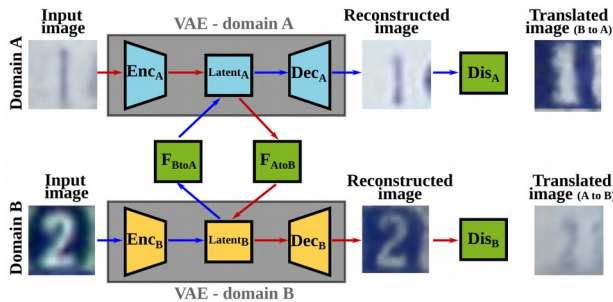
$$\mathcal{L}_{GAN}(Gen_A, Dis_A) = \mathbb{E}_{x \in dom(A)} [\log(Dis_A(x))] \\ + \mathbb{E}_{y \in dom(B)} [\log(1 - Dis_A(Gen_A(y)))]$$

→ ensures the mapping between latent spaces is accurate.

Proposed Method

Shape Translation

- Learn a mapping function between the [shape] domains.
 - Through an extended CycleGAN¹.



→ Applying the Loss

$$\begin{aligned} \mathcal{L} = & \mathcal{L}_{GAN}(Gen_A, Dis_A) + \mathcal{L}_{GAN}(Gen_B, Dis_B) \\ & + \lambda_{cycle} * \mathcal{L}_{cyc}(Gen_A, Gen_B) \\ & + \lambda_{sim} * \mathcal{L}_{sim}(x, y, Gen_A, Gen_B) \end{aligned}$$

Where:

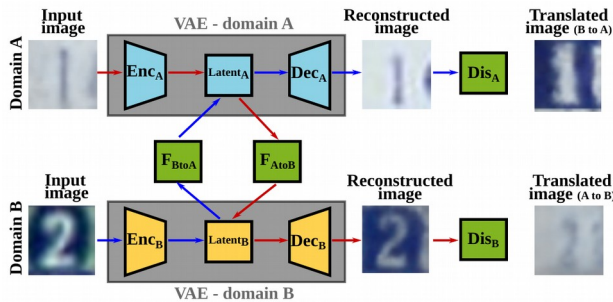
$$\begin{aligned} \mathcal{L}_{cyc} = & \mathbb{E}_{x \in dom(A)} [||Gen_B(Gen_A(x)) - x||_1] \\ & + \mathbb{E}_{y \in dom(B)} [||Gen_A(Gen_B(y)) - y||_1] \end{aligned}$$

→ ensures the cycle consistency to hold.

Proposed Method

Shape Translation

- Learn a mapping function between the [shape] domains.
 - Through an extended CycleGAN¹.



→ Applying the Loss

$$\begin{aligned}\mathcal{L} = & \mathcal{L}_{GAN}(Gen_A, Dis_A) + \mathcal{L}_{GAN}(Gen_B, Dis_B) \\ & + \lambda_{cycle} * \mathcal{L}_{cyc}(Gen_A, Gen_B) \\ & + \lambda_{sim} * \mathcal{L}_{sim}(x, y, Gen_A, Gen_B)\end{aligned}$$

Where:

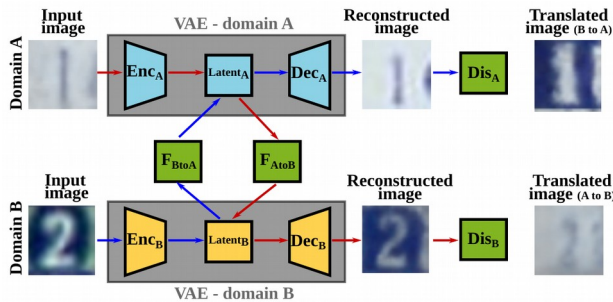
$$\begin{aligned}\mathcal{L}_{sim}(x, y, Gen_A, Gen_B) = & \mathbb{E}_{x \in dom(B)} d(x, Gen_A(x)) \\ & + \mathbb{E}_{y \in dom(A)} d(y, Gen_B(y))\end{aligned}$$

→ favours “good” translations.

Proposed Method

Shape Translation

- Learn a mapping function between the [shape] domains.
 - Through an extended CycleGAN¹.



→ Applying the Loss

$$\begin{aligned} \mathcal{L} = & \mathcal{L}_{GAN}(Gen_A, Dis_A) + \mathcal{L}_{GAN}(Gen_B, Dis_B) \\ & + \lambda_{cycle} * \mathcal{L}_{cyc}(Gen_A, Gen_B) \\ & + \lambda_{sim} * \mathcal{L}_{sim}(x, y, Gen_A, Gen_B) \end{aligned}$$

Where:

$$\mathcal{L}^{SSIM}(x, y) = \frac{1}{N} \sum_{p=1}^N (1 - SSIM(x_p, y_p)) \leftarrow \text{perceptual similarity}^2$$

→ favours “good” translations.

¹Zhu, et al. “Unpaired image- to-image translation using cycle-consistent adversarial networks”, ICCV’17.

²Wang et al., “Multiscale structural similarity for image quality assessment,” Conf.on Signals, Systems Computers, 2003.

Evaluation

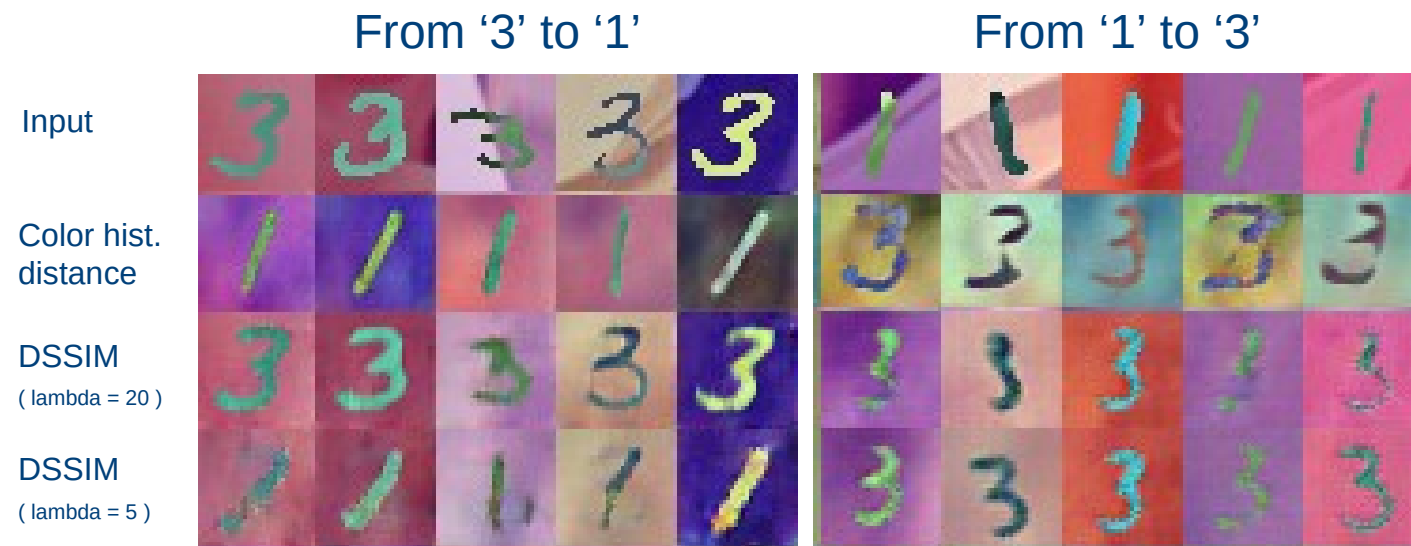
Evaluation

Translating digits in synthetic images

Color MNIST dataset¹

~ 6.5K images per digit class.

- Results



Observations

- The Color loss fails at preserving the style
- The DSSIM loss tends to produce blurry results for $\lambda=20$.

¹<https://www.wouterbulten.nl/blog/tech/getting-started-with-gans-2-colorful-mnist/>

Evaluation

Translating digits in real images

SVHN dataset¹

~ 7.3K images per digit class.

- Qualitative Results



Observations

- Overall translation is good.
- For some classes, translation is slightly blurry. (e.g. class '0' & '6')

¹Coates et al., "Reading digits in natural images with unsupervised feature learning", NIPS'11 Workshops.

Evaluation

Translating digits in real images

SVHN dataset¹

~ 7.3K images per digit class.

- Quantitative Results

Use classification performance as a proxy metric

target class	(a)	(b)	(c)	(d)
0	96.39	85.09	65.95	75.68
1	97.94	98.47	97.71	92.50
2	95.61	95.64	92.74	95.81
3	92.68	95.91	93.29	95.40
4	96.23	96.04	93.80	91.97
5	94.80	89.77	79.58	84.62
6	95.70	88.67	74.25	81.49
7	94.06	94.95	84.37	89.47
8	93.07	86.57	73.41	82.28
9	94.98	89.34	75.09	85.06
Avg.	95.46	93.47	86.51	89.27

Classification performance from:

- a) Original SVHN test set.
- b) Images reconstructed by the VAEs.
- c) Translated Images from class '1'.
- d) Translated Images from class '2'.

¹Coates et al., "Reading digits in natural images with unsupervised feature learning", NIPS'11 Workshops.

Evaluation

Translating digits in real images

SVHN dataset¹

~ 7.3K images per digit class.

- Quantitative Results

Use classification performance as a proxy metric

target class	(a)	(b)	(c)	(d)
0	96.39	85.09	65.95	75.68
1	97.94	98.47	97.71	92.50
2	95.61	95.64	92.74	95.81
3	92.68	95.91	93.29	95.40
4	96.23	96.04	93.80	91.97
5	94.80	89.77	79.58	84.62
6	95.70	88.67	74.25	81.49
7	94.06	94.95	84.37	89.47
8	93.07	86.57	73.41	82.28
9	94.98	89.34	75.09	85.06
Avg.	95.46	93.47	86.51	89.27

Classification performance from:

- a) Original SVHN test set.
- b) Images reconstructed by the VAEs.
- c) Translated Images from class '1'.
- d) Translated Images from class '2'.

Observations

- Performance is not uniform over the classes.

¹Coates et al., "Reading digits in natural images with unsupervised feature learning", NIPS'11 Workshops.

Evaluation

Translating digits in real images

SVHN dataset¹

~ 7.3K images per digit class.

- Quantitative Results

Use classification performance as a proxy metric

target class	(a)	(b)	(c)	(d)
0	96.39	85.09	65.95	75.68
1	97.94	98.47	97.71	92.50
2	95.61	95.64	92.74	95.81
3	92.68	95.91	93.29	95.40
4	96.23	96.04	93.80	91.97
5	94.80	89.77	79.58	84.62
6	95.70	88.67	74.25	81.49
7	94.06	94.95	84.37	89.47
8	93.07	86.57	73.41	82.28
9	94.98	89.34	75.09	85.06
Avg.	95.46	93.47	86.51	89.27

Classification performance from:

- a) Original SVHN test set.
- b) Images reconstructed by the VAEs.
- c) Translated Images from class '1'.
- d) Translated Images from class '2'.

Observations

- Performance is not uniform over the classes.
- In some cases it is as good as on the original images.

¹Coates et al., "Reading digits in natural images with unsupervised feature learning", NIPS'11 Workshops.

Evaluation

Translating digits in real images

SVHN dataset¹

~ 7.3K images per digit class.

- Quantitative Results

Use classification performance as a proxy metric

target class	(a)	(b)	(c)	(d)
0	96.39	85.09	65.95	75.68
1	97.94	98.47	97.71	92.50
2	95.61	95.64	92.74	95.81
3	92.68	95.91	93.29	95.40
4	96.23	96.04	93.80	91.97
5	94.80	89.77	79.58	84.62
6	95.70	88.67	74.25	81.49
7	94.06	94.95	84.37	89.47
8	93.07	86.57	73.41	82.28
9	94.98	89.34	75.09	85.06
Avg.	95.46	93.47	86.51	89.27

Classification performance from:

- a) Original SVHN test set.
- b) Images reconstructed by the VAEs.
- c) Translated Images from class '1'.
- d) Translated Images from class '2'.

Observations

- Performance is not uniform over the classes.
- In some cases it is as good as on the original images.
- Low translations performance seems to be correlated with weak domain modelling.

¹Coates et al., "Reading digits in natural images with unsupervised feature learning", NIPS'11 Workshops.

Evaluation

Translating digits in real images

SVHN dataset¹

~ 7.3K images per digit class.

- Quantitative Results

Use classification performance as a proxy metric

target class	(a)	(b)	(c)	(d)
0	96.39	85.09	65.95	75.68
1	97.94	98.47	97.71	92.50
2	95.61	95.64	92.74	95.81
3	92.68	95.91	93.29	95.40
4	96.23	96.04	93.80	91.97
5	94.80	89.77	79.58	84.62
6	95.70	88.67	74.25	81.49
7	94.06	94.95	84.37	89.47
8	93.07	86.57	73.41	82.28
9	94.98	89.34	75.09	85.06
Avg.	95.46	93.47	86.51	89.27

Classification performance from:

- a) Original SVHN test set.
- b) Images reconstructed by the VAEs.
- c) Translated Images from class '1'.
- d) Translated Images from class '2'.

Observations

- Performance is not uniform over the classes.
- In some cases it is as good as on the original images.
- Low translations performance seems to be correlated with weak domain modelling.

¹Coates et al., "Reading digits in natural images with unsupervised feature learning", NIPS'11 Workshops.

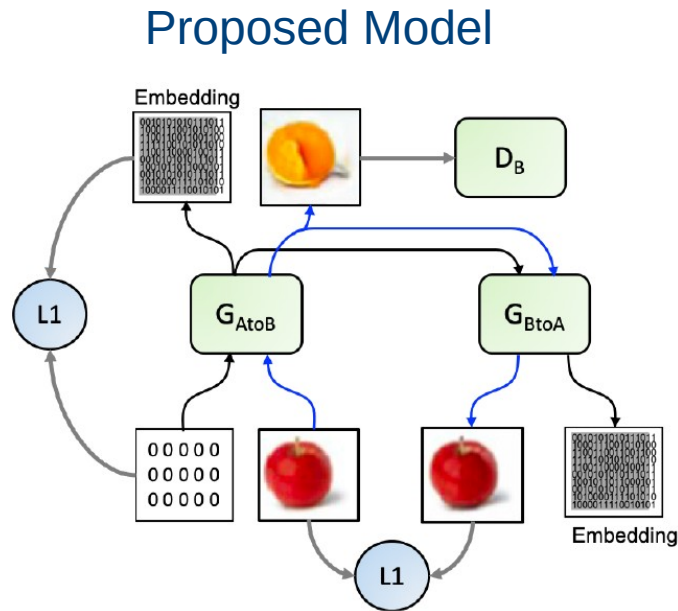
Take-home Message

- **Transferring the structure while preserving the style is possible** → within simple scenarios.
- **Background clutter poses a challenge on this type of translation** → further experimentation is required.

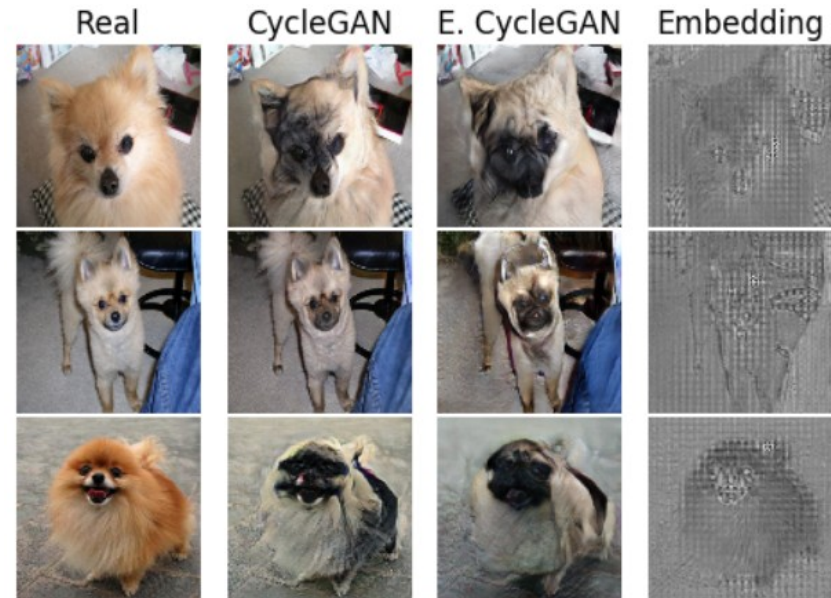
There is Hope

Here at ICIP'19 ! ...

Embedded CycleGAN for Shape-agnostic Translation¹



Some Results



¹ Longman & Ptucha., "Embedded cyclegan for shape-agnostic image-to-image translation", ICIP'19.

Follow-up Work ...

Unpaired Shape Translation¹

Clothing translation



Face translation



¹ Wang et al., "Unsupervised shape transformer for image translation and cross-domain retrieval". ArXiv:12812.02134 .

Follow-up Work ...

Unpaired Shape Translation¹

Clothing translation



Face translation



Cross-domain →
image retrieval



¹ Wang et al., "Unsupervised shape transformer for image translation and cross-domain retrieval". arXiv:12812.02134 .

Contact

José Oramas M.

Email: jose.oramas@esat.kuleuven.be

Website: <http://homes.esat.kuleuven.be/~joramas>

Twitter: [@jaom7](https://twitter.com/jaom7)



Questions?



Towards Object Shape Translation Through Unsupervised Generative Deep Models

Lies Bollens, Tinne Tuytelaars & José Oramas M.

KU Leuven, ESAT-PSI

September 25th, 2019

