

# Synthetic Data and Generative Models in Deep Learning

Sergey Nikolenko

October 13, 2018



# Major limiting factors for AI adoption

Data is usually the most important limiting factor

Current obstacles to AI adoption in various domains:

- lack of sufficiently large datasets
- lack of labeled training data
- difficulty explaining results
- difficulty generalizing
- risk of bias

Here at Neuromation, we are addressing these issues with...

# MCKINSEY GLOBAL INSTITUTE NOTES FROM THE FRONTIER

#### Thinking on its own: Al in the NHS

Clobal Artificial Intelligence (AI) Market Report 2018: Annual Estimates & Forecasts 2016-2024 - Growth Drivers, Trends & Issues

RESEARCHANDMARKETS

#### **GLOBAL ENTERPRISE AI MARKET 2018-2022**





# **Synthetic Data** For Deep Learning in Computer Vision and Beyond



- For computer vision, we generate 100% accurate labeled data, with pixel-perfect labeling which is very hard or impossible to do by hand
- Increases the speed of automation by orders of magnitude and is several times cheaper than hand labeling

# Synthetic data is a possible solution in many cases

- It can be traced to 2000s, but before 2016 synthetic data was not popular
- A few works, mostly related to self-driving cars where datasets are clearly unavailable
- E.g., (Marin et al., 2010) use virtual world training sets for pedestrian detection with classical CV (SVM on HOG features)



UROMATION



VIRTUAL-WORLD LEARNING





- Around 2016 interest spiked
- Now one can find various synthetic datasets for computer vision:
  - SunCG (Song et al., 2017): 45K indoor 3D scenes ready to visualize





- Around 2016 interest spiked
- Compare with Matterport3D (Chang et al., 2017): largest real dataset with panoramic views, 90 buildings



**Textured 3D Mesh** 

NEUROMATION







**Object Instances** 





- Synthetic datasets with city scenes:
  - (Johnson-Robertson et al., 2017):
    Driving in the Matrix synthetic
    Cityscapes
  - (Gaidon et al., 2016): virtual worlds for object tracking — synthetic KITTI

















- Synthetic datasets with specific objects:
  - (Varol et al., 2017): SURREAL
    (Synthetic hUmans for REAL tasks)



**NEUROMATION** 



PART SEGMENT



- These synthetic datasets can be wrapped in synthetic enablers:
  - MINOS (Savva et al., 2017) is a simulator for navigation in synthetic 3D environments







- These synthetic datasets can be wrapped in synthetic enablers:
  - House3D (Wu et al., 2018) is a similar simulator from Facebook
  - And then such simulators can be used for...

NEUROMATION





- Reinforcement learning:
  - RoomNav (Wu et al., 2018) uses House3D to learn to navigate rooms





- Reinforcement learning:
  - (Das et al., 2018) uses House3D for embodied question answering
  - Planner and controller to improve gradient flows



UKRAINE



NEUROMATION



- NLP for navigation:
  - (Anderson et al., 2018) uses Matterport3D to create R2R, a dataset of room-to-room navigation instructions
  - And an action prediction model with attention over the instructions that actually can navigate



Pass the pool and go indoors using the double glass doors. Pass the large table with chairs and turn left and wait by the wine bottles that have grapes by them.

Walk straight through the room and exit out the door on the left. Keep going past the large table and turn left. Walk down the hallway and stop when you reach the 2 entry ways. One in front of you and one to your right. The bar area is to your left.

Enter house through double doors, continue straight across dining room, turn left into bar and stop on the circle on the ground.



Standing in front of the family picture, turn left and walk straight through the bathroom past the tub and mirrors. Go through the doorway and stop when the door to the bathroom is on your right and the door to the closet is to your left.

Walk with the family photo on your right. Continue straight into the bathroom. Walk past the bathtub. Stop in the hall between the bathroom and toilet doorways.

Walk straight passed bathtub and stop with closet on the left and toilet on the right.



# Neuromation Project: The Retail Industry

# **170,000 OBJECTS** TO BE RECOGNIZED ON SHELVES

- Impossible to label this dataset by hand
- We generate synthetic data to train DL models
- On the right: our sample synthetic images.
- Generated by rendering of 3D scenes.

# In computer vision, synthetic data can come from 3d modeling





## Computer Vision For Baby Monitors Healthcare example: synthetic babies to train smart cameras

- **Example from healthcare:** we are collaborating with MonBaby, a baby monitor that helps prevent SIDS and other emergencies through automated alarms
- Current Monbaby solution: a wearable button
- Our joint solution: a smart camera that watches the baby and recognizes the pose, actions, breathing, etc.
- Synthetic data solves the lack of real data, especially for emergencies







## Synthetic Data for Medical Imaging? Success despite lack of data, but how can we go synthetic?..

- Interesting question: how do we apply synthetic data to, say, medical imaging?
- We have a number of projects in medical imaging, showing that they can succeed despite lack of data
- But can we still produce useful synthetic data for domains where direct 3D modeling is hard or impossible?
- What kind of models would it take?..







- ...why, generative models, of course!
- Sometimes we cannot produce synthetic data directly, but can still augment and enlarge datasets with models such as GANs:
  - generator generates
  - discriminator discriminates
  - lots and lots of variations





• Let us look at the story of one specific application...



- (Zhang et al., 2015): gaze estimation in the wild; MPIIGaze dataset, standard convolutional architecture
- Recent state of the art

NEUROMATION













 (Wood et al., 2016): learning gaze estimation from synthetic images generated by a special 3D modeling system UnityEyes

ROMATION





**Figure 8:** We include eyelashes and eye wetness for realism. (a) shows a render without these, (b) shows eye wetness (red) and eyelash geometry (blue), (c) shows the final render.



**Figure 9:** We use pre-integrated skin shading for realism (c). Without it, skin appears too hard (a). (b) shows the scattered light through skin – this causes the skin to appear soft.



 (Shrivastava et al., 2017): Apple learns a gaze estimation model on synthetic images refined by SimGAN that trains to fool D while self-regularizing to keep refined similar to original synthetic images

ROMATION



Unlabeled Real Images





• (Shrivastava et al., 2017): as a result, Apple significantly improved state of the art in gaze estimation; also applied to hand pose estimation

Method	R/S	Error
Support Vector Regression (SVR) [33]	R	16.5
Adaptive Linear Regression ALR) [23]	R	16.4
Random Forest (RF) [36]	R	15.4
kNN with UT Multiview [47]	R	16.2
CNN with UT Multiview [47]	R	13.9
k-NN with UnityEyes [43]	S	9.9
CNN with UnityEyes Synthetic Images	S	11.2
CNN with UnityEyes Refined Images	S	7.8





- And indeed, this idea has been applied to medical imaging
- (Shin et al., Sep 2018): synthetizing MRI images with brain tumors with GANs

Method		Real	Real + Synthetic	Synthetic only	Synthetic only, fine-tune on 10% real
GAN-based (no aug)	0.	64/0.14	0.80/0.07	0.25/0.14	0.80/0.18
GAN-based (with aug)	0.	81/0.13	0.82/0.08	0.44/0.16	0.81/0.09
Wang et al. [20]	0.	85/0.15	0.86/0.09	0.66/0.13	0.84/0.15







BAGAN (balancing GAN) solves another problem: data augmentation for GAN Generator instances highly imbalanced Autoencoder vecto Class-Encoder Decode conditional latent vector datasets vecto Origina generator data Origina econstructed data data In loss optimization (Mariani et al., 2018) instance (a) Autoencoder training.



(b) GAN initialization.



<sup>(</sup>c) GAN training.



(a) Real image samples



(b) BAGAN







(d) Simple GAN





# Synthetic Data Becomes Ubiquitous

- Recent major conferences keep the trend:
  - CVPR 2018: coping with domain shift to train on synthetic data
  - ICLR 2018: new multimodal synthetic datasets
  - ...and much more!

UROMATION







- Somethimes synthetic data is needed directly, not as a stepping stone for training models
- Our collaboration with Insilico Medicine: GANs for drug discovery that generate (fingerprints of) molecules likely to have desired properties with conditional adversarial autoencoders



# AI Talent + Synthetic Data + Privacy + Compute = ?

I have made the case for synthetic data, but there are other problems:

- Sometimes data is sensitive, and data providers do not trust public clouds but still need computational power to train models.
- Most important problem: dire need for AI talent; only ~22K deep learning experts in the world, mostly concentrated in a few places (see map on the right).





# **Neuromation Platform**

Ai Talent + Synth. Data + Compute = Neuromation Platform

Democratizing AI for every industry, healthcare first of all



NEUROMATION



## **Our Team**



Yashar Behzadi CEO



**Constantine Goltsev** Investor / Chairman



**Evan Katz** Chief Revenue Officer



Sergey Nikolenko Chief Research Officer



СТО



Maxym Prasolov Founder



Yuri Kundin COO



**Evgeniya Zaslavskaya** PR & Business Development (Russia & CIS)



Arthur McCallum VP Digital Economy



Fedor Savchenko VP of Research & Development



**David Orban** Adviser



Andrew Rabinovich Adviser





# **Our Research Team**



#### Andrew Rabinovich Adviser

World-leading researcher in deep learning and computer vision research, author of numerous patents and publications, founder of a biotech startup; held leading R&D positions at Google; currently the director of deep learning at Magic Leap.



#### Sergey Nikolenko Chief Research Officer

World-class researcher in machine learning (deep learning, Bayesian methods, NLP, and more) and analysis of algorithms (for networking), Sergey has authored more than 120 research papers, several books, courses on machine learning and deep learning, and more.



Kyryl Truskovskyi Lead Researcher

ROMATION



Aleksey Artamonov Lead Researcher





Alexander Rakhlin Senior Researcher Kaggle Master



#### +10 more top deep learning researchers in our offices at St. Petersburg, Kiev, and San Francisco!





# Where androids dream of electric sheep

THANK YOU!







# Breast Cancer Histology

#### Medical imaging projects: success despite lack of data

scale, crop size

**Rakhlin et al., 2018.** Breast cancer histology image analysis with deep CNNs, four classes: normal, benign, *in situ* carcinoma, invasive carcinoma





Preprocessing pipeline





5. 3-norm pooling

Augmented crops



4. Deep CNN descriptors





6. Building multiple datasets, based on number of scales, crop sizes, encoders







#### Pediatric Bone Age Assessment Medical imaging projects: success despite lack of data

**Iglovikov et al., 2017.** Skeletal bone age assessment used to diagnose endocrine and metabolic disorders in child development; deep CNNs for keypoint detection and segmentation









#### Diabetic Retinopathy Detection Medical imaging projects: success despite lack of data

**Rakhlin et al., 2017.** Diabetic retinopathy detection surpassing human optometrists and achieving state of the art results with much less data







# Medical Imaging Beyond Computer Vision

**Deep learning methods are not restricted to 2D/3D images** 

- Imaging mass-spectrometry: spatially structured multidimensional data
- Project led by Dr. Theodore Alexandrov at the European Molecular Biology Lab
- High-def imaging, 10K+ dimensional spectrum at every point
- We are developing informative latent representations to help study the cell cycle etc. via metabolomics

