# AutoBayes:
# Automated Machine Learning with Bayesian Graph Exploration for Nuisance-Robust Inference

Toshiaki Koike-Akino

Andac Demir

Ye Wang

Deniz Erdogmus

Aug. 2020

**AUTOMATED MACHINE LEARNING (AUTOML):**
Breakthrough Technology to Accelerate Machine Learning Outcomes

Thomas Bayes
1702 - 1761

- Part I: Trends of machine learning

- Part II: **Adversarial learning** for nuisance-robust data analysis

- Part III: **Meta learning**: Automated machine learning (AutoML)
  - Automated architecture and hyperparameter tuning

- Part IV: AutoBayes
  - Bayesian inference graph modeling
  - Bayes Ball algorithm
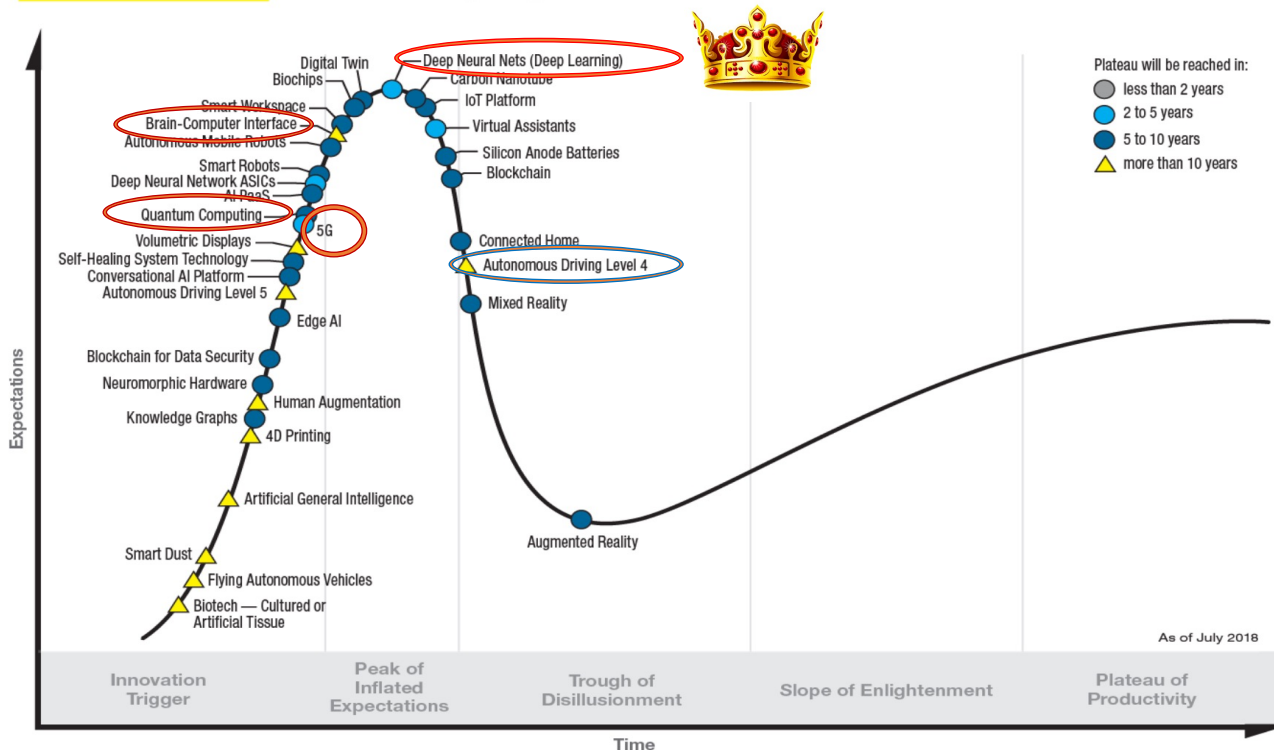
- Part V: **Ensemble learning**

- Summary

**AUTOMATED MACHINE LEARNING (AUTOML):**
Breakthrough Technology to Accelerate Machine Learning Outcomes

*Answering...*
What is best DNN architecture?
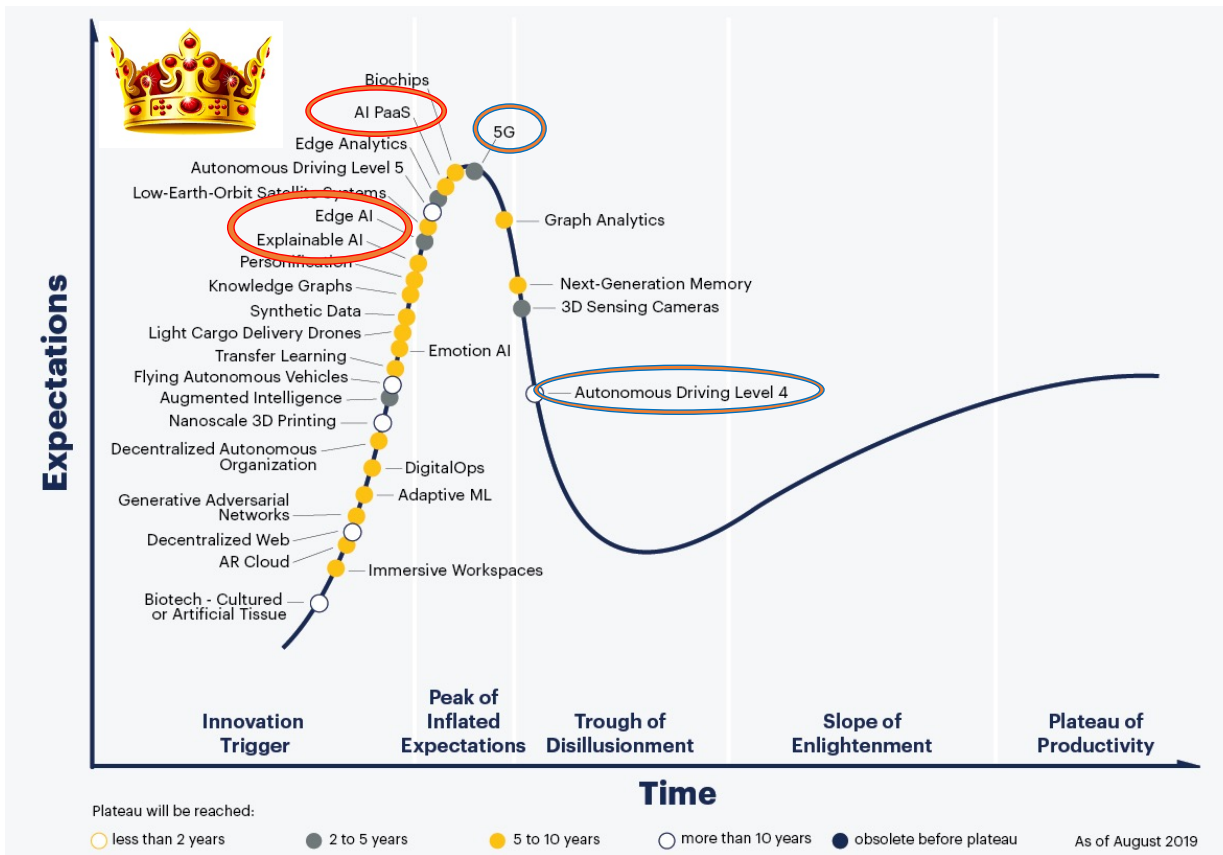
# Emerging Technologies 2018

- Gartnar's Hype Cycle for Emerging Technologies, 2018 July
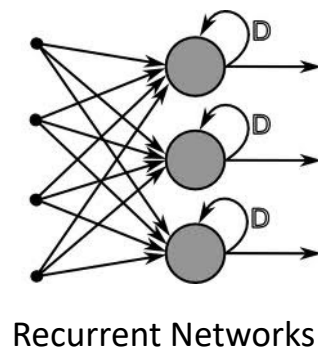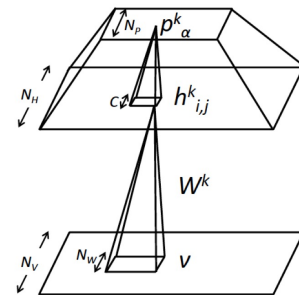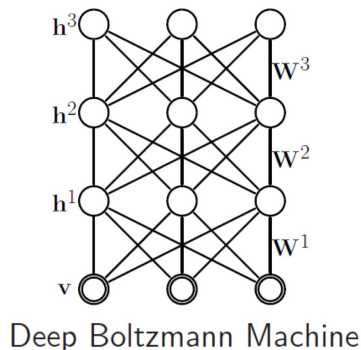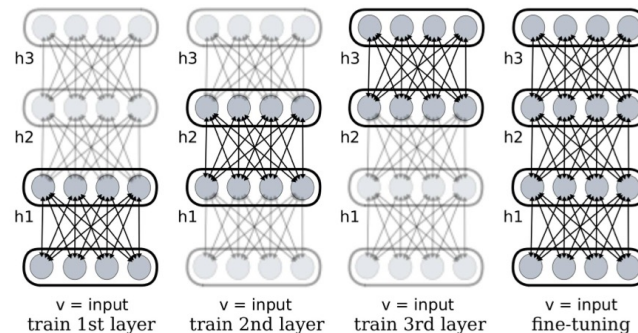
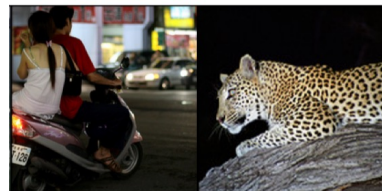- Gartnar's Hype Cycle for Emerging Technologies, 2019 August

# Deep Learning (DL) for Artificial Intelligence (AI)

- Deep learning = fancy name of multi-layer perceptron neural networks.
  - 2006 Hinton: Many layers, layer-wise pre-training, massive data sets

- Massively parallel computation
  - Driver: graphic processor units, tensor processor units …

- Variants:
  - Deep belief networks
  - Deep convolutional networks
  - Deep recurrent networks
  - Deep Boltzmann machines
  - Deep autoencoder



v = input
train 1st layer

v = input
train 2nd layer

v = input
train 3rd layer

v = input
fine-tuning



Deep Boltzmann Machine

Deep Belief Network

Convolutional Networks

Recurrent Networks

# Deep Learning for Media Signal Processing

- Audio & Visual Applications



motor scooter | leopard
"man in black shirt is playing guitar."

# AI Surpassing Human-Level Performance

# Moore's Law: Exponential Growth in Applications

- Hit count of articles per year in GoogleScholar; *Wireless Communication* applications



Figure showing "Number of Articles" vs "Year" with legend:
- Machine Learning + Wireless Comm
- Deep Learning + Wireless Comm
- Exponential Prediction

Machine Learning — 124% annually — 1.2ˣ

Koike, Neural MIMO detection, WPMC 2004

Deep Learning — 212% annually — 2.1ˣ

# Moore's Law: Exponential Growth in Applications II

- Hit count of articles per year in GoogleScholar; *Optical Communication* applications

# Applied Deep Learning

- AI has been applied to various fields

Wireless Communication

Networked Control

Optical Communication

Localization Navigation

Device / Integrated Circuit

Tomography Imaging
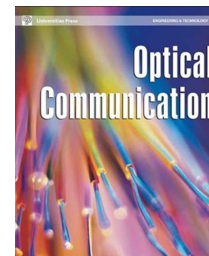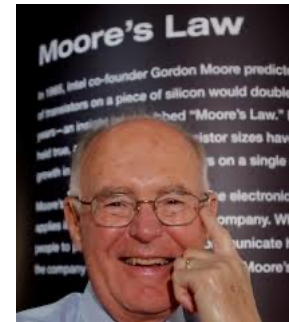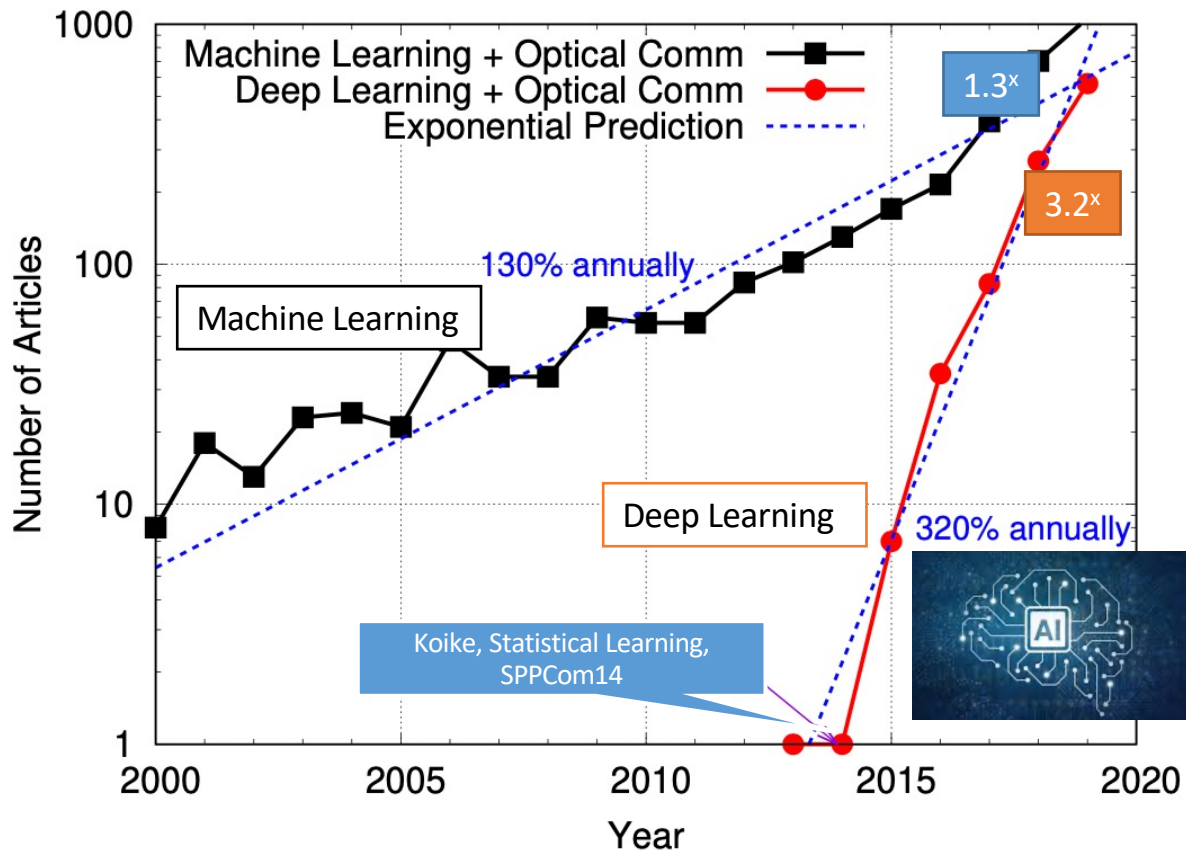
Bio-Sensing Human Interface

AI: Essential Component for R&D

# Biosignal Processing and Mind Sensing

- Joint work with Prof. Deniz Erdogmus (Northeastern Univ.)

- 2015 Ruhi Mahajan
  - Authentication (EMBC)

- 2016 Fernando Quivira
  - Probabilistic GMM+LSTM (BHI)

- 2017 Chun-Shu Wei
  - Few-shot learning (NER)

- 2018 Ozan Ozdenizci
  - Adversarial VAE (NER, SPL, Access)

- 2019 Mo Han
  - Complementary adversarial (EMBC, SPL)
  - Rateless soft disentangling (JBHI)

- 2020 Andac Demir
  - AutoBayes (Access)
  - Graph EEG net (EMBC)

Northeastern Univ.

# Our Publications

1. Koike-Akino, T., Mahajan, R., Marks, T.K., Tuzel, C.O., Wang, Y., Watanabe, S., Orlik, P.V., **"High-Accuracy User Identification Using EEG Biometrics"**, IEEE EMBC, August 2016.

2. Wang, Y., Koike-Akino, T., Erdogmus, D. "**Invariant Representations from Adversarially Censored Autoencoders**", arxiv:1805.08097, May 2018.

3. Quivira, F., Koike-Akino, T., Wang, Y., Erdogmus, D., **"Translating sEMG Signals to Continuous Hand Poses using Recurrent Neural Networks"**, IEEE BHI, January 2018.

4. Wei, C.-S., Koike-Akino, T., Wang, Y., **"Spatial Component-wise Convolutional Network (SCCNet) for Motor-Imagery EEG Classification"**, IEEE NER, March 2019.

5. Ozdenizci, O., Wang, Y., Koike-Akino, T., Erdogmus, D., **"Transfer Learning in Brain-Computer Interfaces with Adversarial Variational Autoencoders"**, IEEE NER, March 2019. (arxiv:1812.06857, Dec. 2018)
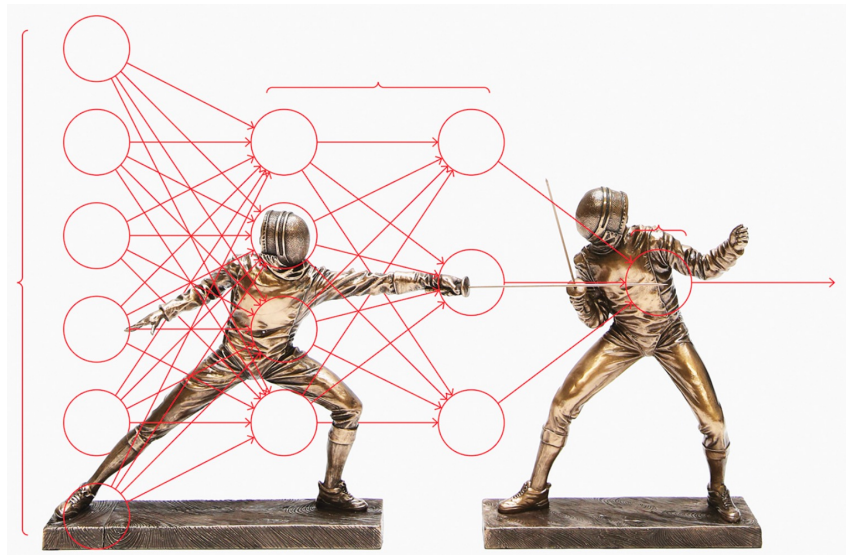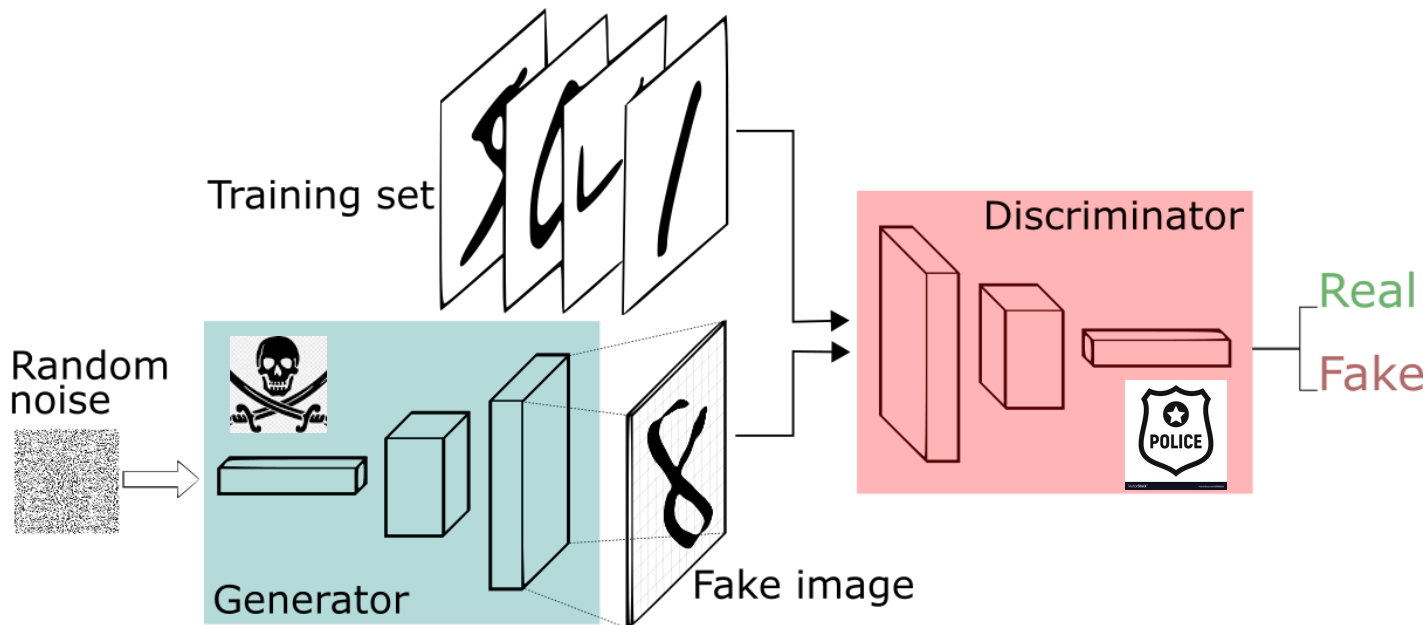
6. Ozdenizci, O., Wang, Y., Koike-Akino, T., Erdogmus, D., **"Adversarial Deep Learning in EEG Biometrics"**, IEEE SPL, March 2019. (arxiv:1903.11673, May 2019)

7. Ozdenizci, O., Wang, Y., Koike-Akino, T., Erdogmus, D., **"Learning Invariant Representations from EEG via Adversarial Inference"**, IEEE Access , April 2020.

8. Koike-Akino, T., Wang, Y., **"Stochastic Bottleneck: Rateless Auto-Encoder for Flexible Dimensionality Reduction"**, IEEE ISIT, June 2020. (arxiv:2005.02870, May 2020)

9. Han, M., Ozdenizci, O., Wang, Y., Koike-Akino, T., Erdogmus, D., **"Disentangled Adversarial Transfer Learning for Physiological Biosignals"**, IEEE EMBC, July 2020. (arxiv:2004.08289, Apr. 2020)

10. Han, M., Ozdenizci, O., Wang, Y., Koike-Akino, T., Erdogmus, D., **" Disentangled Adversarial Autoencoder for Subject-Invariant Physiological Feature Extraction "**, IEEE SPL, July 2020. (arxiv:2008.11426)

11. Demir, A., Koike-Akino, T., Wang, Y., Erdogmus, D., **"AutoBayes: Automated Bayesian Graph Exploration for Nuisance-Robust Inference"**, IEEE Access, Mar. 2021. (*arxiv:2007.01255*, July 2020).

12. Haruna, M., Ogino, M., Koike-Akino, T., "**Proposal and Evaluation of Visual Haptics for Manipulation of Remote Machine System,**" Frontiers, Aug. 2020.

13. Han, M., Ozdenizci, O., Wang, Y., Koike-Akino, T., Erdogmus, D., **" Universal Physiological Representation Learning with Soft-Disentangled Rateless Autoencoders"**, IEEE JBHI, Mar. 2021. arxiv:2009.13453

# Adversarial Learning

# Adversarial Networks

- Generative Adversarial Networks (GAN) [Goodfellow et al, 2014]
  - Train two **competing** neural networks
  - Generator learns to fake images by trying to fool Discriminator



- Competition between counterfeiters and police ⇒ better fake money

Aug. 2020                                                                                                    16

# GAN for Synthetic Faces

- Nvidia GAN Results [Karras et al, 2018]



Realistic Fake Faces
youtube:XOxxPcy5Gr4
youtube:kSLJriaOumA

# CycleGAN for Image Translation

- CycleGAN [Zhu et al, 2017]



© MERL       Aug. 2020

Style Change

- Many different ways to adversarially combine networks



(a) Encoder network E

(b) Noiseless joint PPGN-h

(c) Joint PPGN-h

# Learning Nuisance-Invariant Data Representations



- Objective: extract invariant representations (features)
  - Remove nuisance variations, sensitive attributes
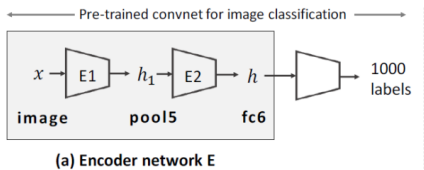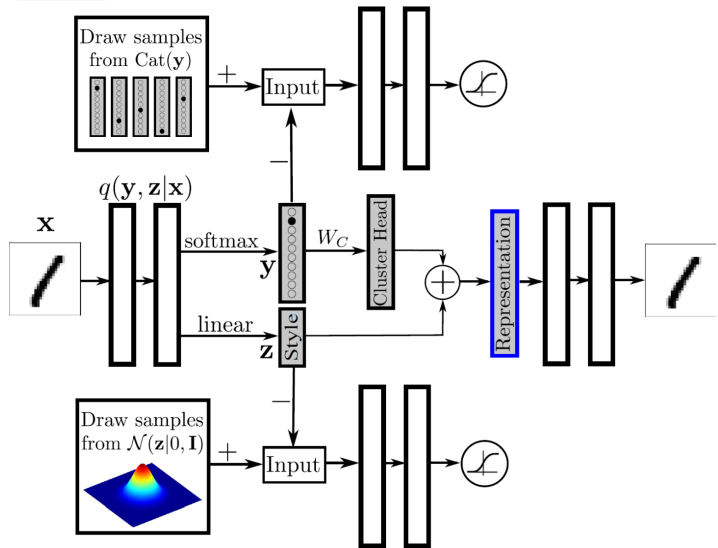  - Motivation: transferability, generalizability, robustness, privacy, fairness
- Autoencoder model: data $\mathbf{x} \rightarrow \boxed{\text{Enc}} \rightarrow$ latent $\mathbf{z} \rightarrow \boxed{\text{Dec}} \rightarrow \hat{\mathbf{x}}$
  - General purpose data representations $\mathbf{z}$
  - Can also support translation, feature/style transfer

# Variational AutoEncoders (VAE)

- VAE introduced by [Kingma & Welling, 2014] with conditional extension [Sohn et al, 2015]

- Learn CVAE model: $(\mathbf{x}, \mathbf{s}, \mathbf{z}) \sim p_\theta(\mathbf{x}|\mathbf{z}, \mathbf{s})p(\mathbf{s})p(\mathbf{z})$
  - $\mathbf{x}$ raw data features
  - $\mathbf{s}$ nuisance variations (conditioning variable)
  - $\mathbf{z}$ latent (unobserved) representation
  - **Invariance:** model explicitly specifies independence between $\mathbf{s}$ and $\mathbf{z}$
  - Generative model $p(\mathbf{x}|\mathbf{z}, \mathbf{s})$ from appropriate parametric family
  - Convenient latent model $p(\mathbf{z}) = N(\mathbf{0}; \mathbf{I})$
  - Nuisance model $p(\mathbf{s})$ arbitrary (not used for training)



Writer: S

Encoder          Decoder

# VAE Training



Label $\mathbf{s}$

Data $\mathbf{x}$ → | Encoder $q_\phi(\mathbf{z}|\mathbf{x},\mathbf{s})$ | Latent $\mathbf{z}$ | Decoder $p_\theta(\mathbf{x}|\mathbf{z},\mathbf{s})$ | $\max_{\phi,\theta}$ ← $\mathrm{E}[\log p_\theta(\mathbf{x}|\mathbf{z},\mathbf{s})]$

$\max_\phi$

$-\mathrm{KL}(q_\phi(\mathbf{z}|\mathbf{x},\mathbf{s})||p(\mathbf{z}))$

- Decoder: generative model $p_\theta(\mathbf{x}|\mathbf{z},\mathbf{s})$
- Encoder: variational posterior $q_\phi(\mathbf{z}|\mathbf{x},\mathbf{s})$
  - In principle, $q_\phi(\mathbf{z}|\mathbf{x},\mathbf{s}) \to p_\theta(\mathbf{z}|\mathbf{x},\mathbf{s})$ and hence $\mathbf{z} \perp\!\!\!\perp \mathbf{s}$
- However, in practice, invariance $(I(\mathbf{s};\mathbf{z}) = 0)$ needs to be enforced

$$\max_{\theta,\phi} \mathcal{L}(\theta,\phi) - \lambda I(\mathbf{s};\mathbf{z})$$
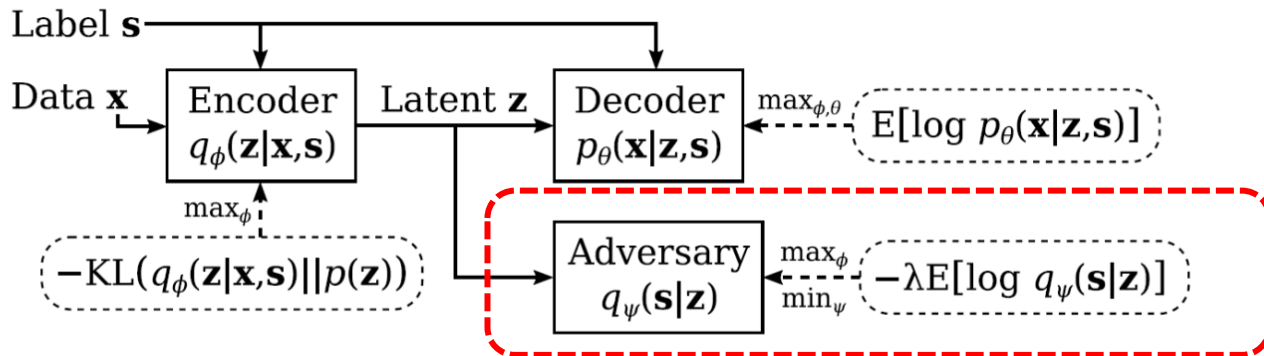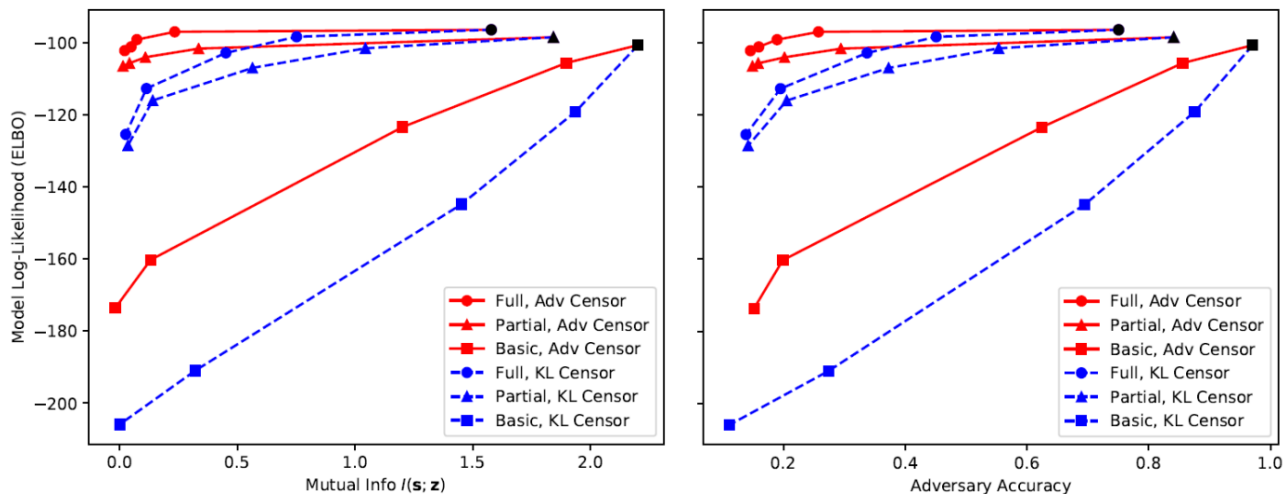
# VAE Training with Adversarial Censoring



- Decoder: generative model $p_\theta(\mathbf{x}|\mathbf{z}, \mathbf{s})$
- Encoder: variational posterior $q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{s})$
  - In principle, $q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{s}) \to p_\theta(\mathbf{z}|\mathbf{x}, \mathbf{s})$ and hence $\mathbf{z} \perp\!\!\!\perp \mathbf{s}$
- However, in practice, invariance $(I(\mathbf{s}; \mathbf{z}) = 0)$ needs to be enforced

$$\max_{\theta, \phi} \mathcal{L}(\theta, \phi) - \lambda I(\mathbf{s}; \mathbf{z}) \Rightarrow \max_{\theta, \phi} \min_\psi \mathcal{L}(\theta, \phi) \underbrace{-\lambda \mathbb{E}\left[\log q_\psi(\mathbf{s}|\mathbf{z})\right]}_{\geq -\lambda(I(\mathbf{z}; \mathbf{s}) - h(\mathbf{s}))}$$

- Adversary $q_\psi(\mathbf{s}|\mathbf{z})$ attempts to recover $\mathbf{s}$, approximates $I(\mathbf{s}; \mathbf{z})$

- Wang, Y., Koike-Akino, T., Erdogmus, D. "**Invariant Representations from Adversarially Censored Autoencoders**", arxiv:1805.08097, May 2018.
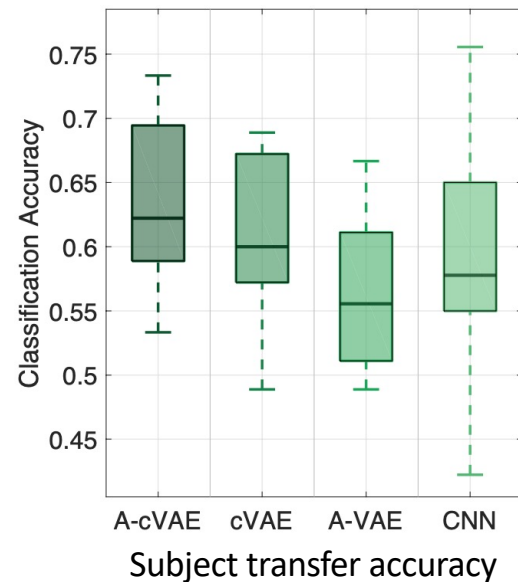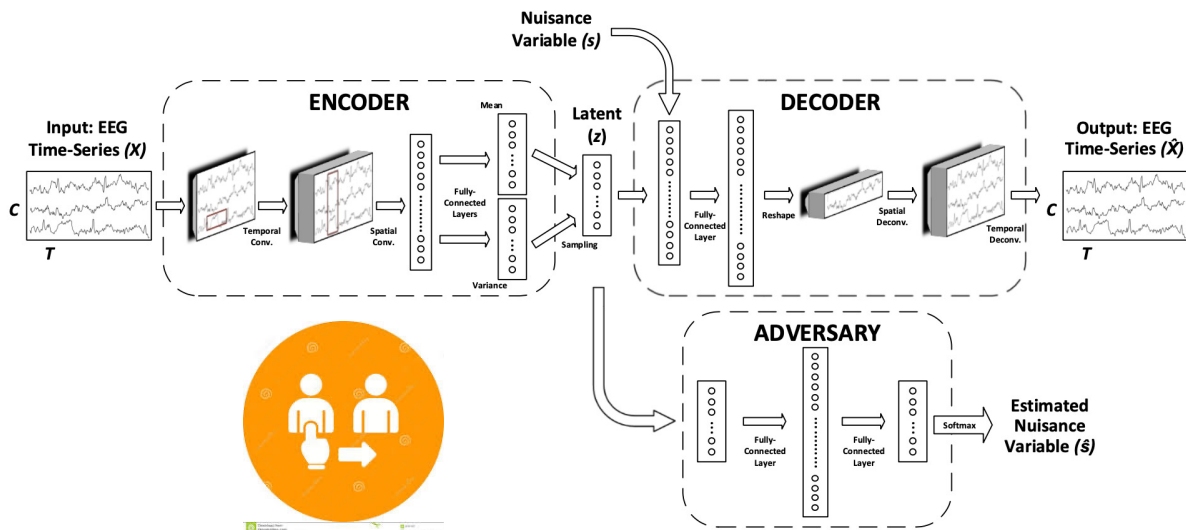


- Full (●): full VAE conditioned on $\mathbf{s}$, i.e. $E(\mathbf{x}, \mathbf{s})$, $D(\mathbf{z}, \mathbf{s})$
- Partial (▲): only decoder conditioned on $\mathbf{s}$, i.e. $E(\mathbf{x})$, $D(\mathbf{z}, \mathbf{s})$
- Basic (■): no conditioning on $\mathbf{s}$, i.e. $E(\mathbf{x})$, $D(\mathbf{z})$
- KL censoring alternative: use $-\gamma \mathrm{KL}\big(q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{s}) \big\| p(\mathbf{z})\big)$, with $\gamma > 1$
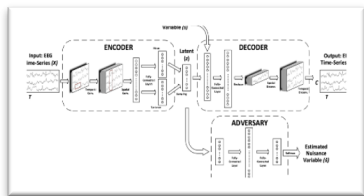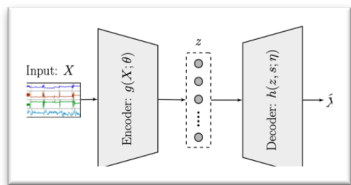
# A-CVAE for Nuisance-Robust Transfer Learning: Zero-Shot Learning

- Cross-subject transfer learning for BCI [Ozdenizci et al, NER'19]
  - Task: motor-imagery decoding from EEG measurements
  - Subject variability is the nuisance variation suppressed

- Cross-session EEG-based biometrics [Ozdenizci et al, SPL'19]
  - Task: subject identification from EEG measurements
  - Session variability is the nuisance variation suppressed





Subject transfer accuracy

# Evolution Map: Nuisance-Invariant Feature Extraction

Rateless soft disentangling

Complementary disentangling

AutoBayes

DA-cRAE

DA-cAE

A-cAE
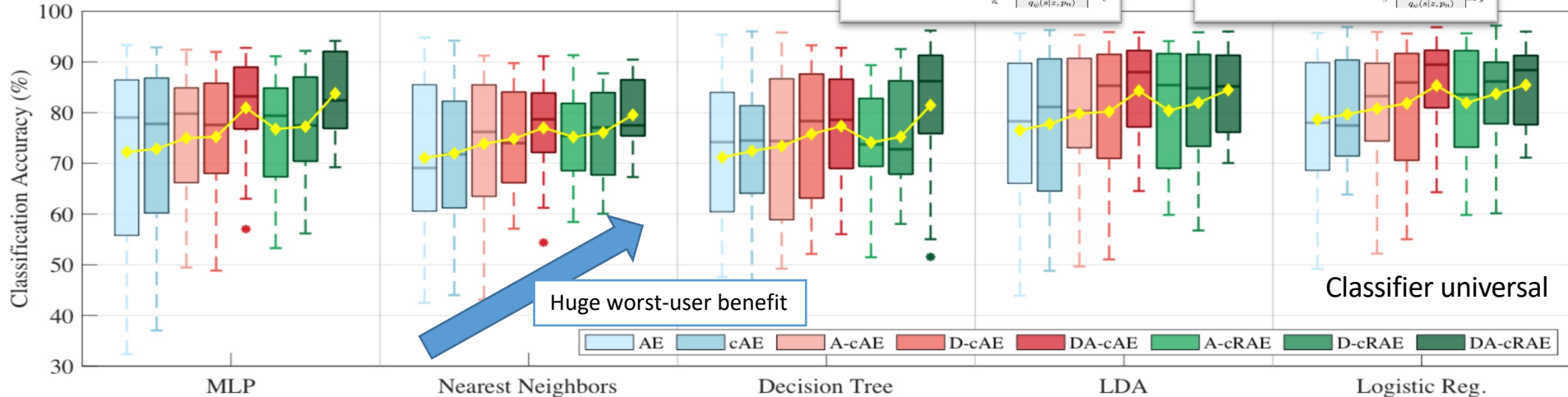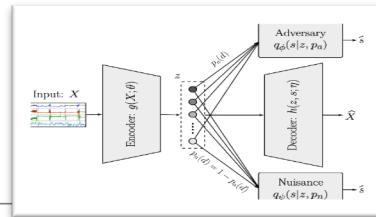
cAE

AE

Classifier universal

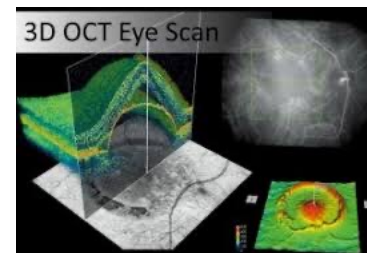Huge worst-user benefit

Classification Accuracy (%)

AE | cAE | A-cAE | D-cAE | DA-cAE | A-cRAE | D-cRAE | DA-cRAE

MLP    Nearest Neighbors    Decision Tree    LDA    Logistic Reg.

Subject transfer accuracy: Stress dataset (heartrate, temperature, etc.)

# Nuisance-Robust Analysis: Not Only for Biosignal Applications

- Nano-photonic device design

- Privacy preserving

- Localization

- Image sensing

- Speech recognition

- Face recognition

- ...

**AutoML**

# Various DNN Architectures

- Forward, recursive, convolutional

- LSTM, GRU, Transformer

- Bottleneck, U-net, HR-net

- ResNet, loopy, clique

- Inception, bilinear



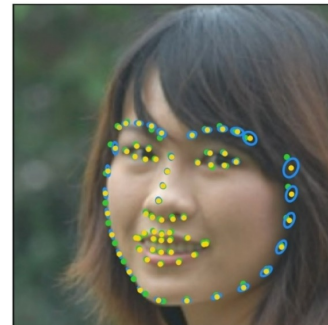End to end to mapping from EBSD patterns to crystallographic orientations



Figure 1: The Transformer - model architecture.

Figure 2. Residual learning: a building block.

MITSUBISHI ELECTRIC
Changes for the Better

*"**Undergraduate-student gradient**" trial-and-error*

- Loss functions: CrossEntropy, MSE, MAE, hinge, L1, …

- Updaters: SGD, Adam, AdaDelta, AdaGrad, AdaMax, LBFGS, RMSprop, …

- Learning rate schedulers: CosineAnnealing, Cyclic, Exponential, MultiStep, OnPlateau, …

- Regularizations: Dropout, Batch Norm, Spectral Norm, DropConnect, StochasticDepth, ShakeDrop, Shake-Shake, …

- Activations: ReLU, sigmoid, tanh, …

- Augmentation, depth, width

- Quantization, initialization

- Pooling, unpooling, padding, …



Batch Normalization [5] in training mode.

# No Free Lunch Theorem

- Wolpert 1996: *What Does Dinner Cost?*, NASA Ames Research Center

- There is no one model that works **best for every problem**

- We thus shall **try multiple models** and find one that works **best for a particular problem**

Human Experts
Programming

AI Experts
Evolutionary Programming

```
model = nn.Sequential(
        nn.Conv2d(1,20,5),
        nn.ReLU(),
        nn.Conv2d(20,64,5),
        nn.ReLU()
    )
```

```
def Setup():
    s4 = 0.5
def Predict(v0):
    v1 = v0 - v9
    v5 = v0 + v9
    m1 = s2 * m2
def Learn(v0, s0):
    s4 = s0 - s1
    s3 = abs(s1)
```

parent

child

Type (i)

```
def Setup():
    s4 = 0.5
def Predict(v0):
    v1 = v0 - v9
    v5 = v0 + v9
    m1 = s2 * m2
def Learn(v0, s0):
    s4 = s0 - s1
    s2 = sin(v1)
    s3 = abs(s1)
```

Google AI

AutoML-Zero

# AutoML, Learning to Learn (L2L), Meta Learning

- Hyperparameter exploration: Auto-Pytorch, …
  - Bayesian optimization
  - Evolutionary optimization

- Architecture exploration
  - Reinforcement learning
  - Cell-based building

- Augmentation exploration

# Linking/Unlinking DNN Cells

- Automated neural architecture search may work, but having little justification
- Why linking, pruning, inserting, stacking, branching, merging, …



Any justification?

**AutoBayes**

* Different from NASA's *AutoBayes…*

# Bayesian Graphical Model

- **Directed graph** indicates marginal dependency



- Joint probability factorization:

$$p(y, s, z, x) = p(y)p(s|y)p(z|s, y)p(x|z, s, y)$$

  - **X**: Measurement (Image, EEG, EMG, …)
  - **Y**: Label to classify (digit, mental state, …)
  - **Z**: Latent variables (reduced-dimension feature, …)
  - **S**: Nuisance variations (user, session, environment, …)
  - Nobody knows true data model…

(4! Possible factorization chains)



Data generative model

# Inference Strategy Justified by Bayesian Graph Model

- If true data model follows Markov: $X - Y$

$$p(y)p(s|y)p(z|s,y)p(x|z,s,y)$$

- Likelihood is independent of $S$ and $Z$

$$p(y,s,z|x) = p(z|x)p(s|z,x)p(y|s,z,x)$$

- Simplest classifier model $p(y|x)$ is sufficient
  - *A-CVAE? - No. Irrational to involve more functionality*
  - *I came up with cool complex model. - No, no, no way!*



Label $Y=4$

Image $X$



(a) Standard Classifier Net



(b) Adversarial CVAE-Based Classifier Net



(c) Potentially Extended Classifier Net

# Inference Strategy Justified by Bayesian Graph Model (II)

- Consider latent-involved Markov: $X - Z - Y$

$$p(y)p(s|\cancel{y})p(z|\cancel{s}, y)\textcolor{blue}{p(x|z, \cancel{s}, \cancel{y})}$$

- Then, we have the likelihood

$$p(z|x)p(s|\cancel{z}, \cancel{x})p(y|\cancel{s}, z, \cancel{x})$$

- *I obtained X-Z-Y network. – No. Not enough*
- *VAE: I added generative model from Bayesian graph. – Not yet*
- *CVAE? – No, X is independent of S*
- *A-VAE? – Yes, reasonable choice*



(b) Model B
Bayesian Graph

Inference Graph

Generative model: VAE

Naïve Inference

*Z* is independent of *S*

- Bayesian graph yields justified inference graphs



Bayesian graph

Inference graph

(a) Model A  (b) Model B  (c) Model C  (d) Model D  (e) Model E  (f) Model F

(g) Model G  (h) Model H  (i) Model I  (j) Model J  (k) Model K

$$p(y,s,z|x) = \begin{cases} p(z|x)p(s|z,x)p(y|s,z,x), & \text{Z-first-inference} \\ p(s|x)p(z|s,x)p(y|z,s,x), & \text{S-first-inference} \end{cases}$$

(b) $Z$-First Inference

(c) $S$-First Inference

Redundant links can be identified

# Bayes Ball Algorithm

- Conditional independence can be justified systematically with simple **10 rules**



An undirected path is active if a Bayes ball travelling along it never encounters the "stop" symbol: ⟶⊣

If there are no active paths from $X$ to $Y$ when $\{Z_1, \ldots, Z_k\}$ are shaded, then $X \perp\!\!\!\perp Y \mid \{Z_1, \ldots, Z_k\}$.

Bayes Ball 10 Rules

Example

no active paths
$$X \perp\!\!\!\perp Y \mid Z$$

one active path
$$X \not\!\perp\!\!\!\perp Y \mid \{W, Z\}$$

- Bayes-Ball finds redundant links for inference graphs



Bayesian Graphs

Bayes-Ball

Inference Graphs

(a) Model A (b) Model B (c) Model C (d) Model D (e) Model E (f) Model F
(g) Model G (h) Model H (i) Model I

(a) Model Dz (b) Model Ds (c) Model Ez (d) Model Es (e) Model Fz (f) Model Fs
(g) Model Gz (h) Model Gs (i) Model Jz (j) Model Js (k) Model Kz (l) Model Ks

- Adversarial alternating updates

$$(\theta, \psi, \eta, \mu) = \arg\min_{\theta, \psi, \eta, \mu} \mathbb{E}\left[\mathcal{L}(\hat{y}, y) + \lambda_s \mathcal{L}(\hat{s}, s) + \lambda_x \mathcal{L}(\hat{x}', x) + \lambda_z \mathbb{KL}(z_1, z_2) - \lambda_a \mathcal{L}(\hat{s}', s)\right],$$

$$(z_1, z_2) = p_\theta(x), \quad \hat{y} = p_\psi(z_1, z_2), \quad \hat{s} = p_\phi(z_1), \quad \hat{x}' = p_\mu(z_1), \quad \hat{s}' = p_\eta(z_1, z_2),$$



Bayesian Graph
Model K

Z-Inference Graph

DA-VAE [Han et al. SPL'20]

- How to connect **Encoder**, **Decoder**, **Classifier**, **Estimator**, **Adversary** cells?

---

**Algorithm 1** Pseudocode for AutoBayes Framework

---

**Require:** Nodes set $\mathcal{V} = [Y, X, S_1, S_2, \ldots, S_n, Z_1, Z_2, \ldots, Z_m]$, where $Y$ denotes task labels, $X$ is a measurement data, $\mathcal{S} = [S_1, S_2, \ldots, S_n]$ are (potentially multiple) semi-supervised nuisance variations, and $\mathcal{Z} = [Z_1, Z_2, \ldots, Z_m]$ are (potentially multiple) latent vectors

**Ensure:** Semi-supervised training/validation datasets

1: **for all** permutations of node factorization from $Y$ to $X$ **do**
2:     Let $\mathcal{B}_0$ be the corresponding Bayesian graph for the permuted full-chain factorization $p(y) \cdots p(z_1 | \cdots) \cdots p(x | \cdots)$
3:         **for all** combinations of link pruning on the full-chain Bayesian graph $\mathcal{B}_0$ **do**
4:             Let $\mathcal{B}$ be the corresponding pruned Bayesian graph
5:             Apply the Bayes-Ball algorithm on $\mathcal{B}$ to build a conditional independency list $\mathcal{I}$
6:             **for all** permutations of node factorization from $X$ to $Y$ **do**
7:                 Let $\mathcal{F}_0$ be the corresponding factor graph, representing a full-chain conditional probability $p(\cdot | x) \cdots p(z_1 | \cdots) \cdots p(y | \cdots, x)$
8:                 Prune all redundant links in $\mathcal{F}_0$ based on conditional independency $\mathcal{I}$
9:                 Let $\mathcal{F}$ be the pruned factor graph
10:                 Merge the pruned Bayesian graph $\mathcal{B}$ into the pruned factor graph $\mathcal{F}$
11:                 Attach an adversary network $\mathcal{A}$ to latent nodes $\mathcal{Z}$ for $Z_k \perp \mathcal{S} \in \mathcal{I}$
12:                 Assign an encoder network $\mathcal{E}$ for $p(\mathcal{Z} | \cdots)$ in the merged factor graph
13:                 Assign a decoder network $\mathcal{D}$ for $p(x | \cdots)$ in the merged factor graph
14:                 Assign a nuisance indicator network $\mathcal{N}$ for $p(\mathcal{S} | \cdots)$ in the merged factor graph
15:                 Assign a classifier network $\mathcal{C}$ for $p(y | \cdots)$ in the merged factor graph
16:                 Adversary train the whole DNN structure with variational reparameterization to minimize a loss function in (11)
17:             **end for**                    ▷ At most $(|\mathcal{V}| - 2)!$ combinations
18:         **end for**                    ▷ At most $2^{|\mathcal{V}|(|\mathcal{V}|-1)/2}$ combinations
19: **end for**                    ▷ At most $(|\mathcal{V}| - 2)!$ combinations
20: **return** the best model having highest task accuracy in validation sets

Automatic exploration of Bayesian graphs

Bayes Ball to check independence

Inference model construction

Link Encoder, Decoder, Classifier, Estimator, and Adversary Nets

Return best architectures

# MITSUBISHI ELECTRIC
*Changes for the Better*

# MNIST (QMNIST)

- 28x28 gray-scale images

- 10-class hand-written digits

- 60,000 training data

- 10,000 test data

- Who wrote?

- QMNIST
  https://github.com/facebookresearch/qmnist
  - Identical datasets of MNIST
  - Extended labels (writer ID etc.)
  - Training data were written by **539** NIST employees
  - Testing data were written by **400** high-schoolers

  - **Writer ID** is a nuisance: Hand-written digits may depend on the writer

- Up to **0.5% gain** by nuisance-robust inference

# Public Physiological Datasets

- Stress: temperature, **heart rate**, electrodermal activity, arterial oxygen level, etc. for 4-state stress level measurement

- RSVP: **EEG** for rapid serial visual presentation (RSVP) drowsiness test with 4 tasks

- MI: PhysioNet EEG Motor Imagery (MI) dataset with 4-class tasks

- ErrP: An error-related potential (ErrP) of EEG dataset in spelling task

- Faces: An implanted electrocorticography (**ECoG**) array dataset for visual stimulus.

- Ninapro: An electromyogram (**EMG**) dataset for fingers motion detection for prosthetic hands.

Heart rates

Electrodermal

Oxygen

RSVP EEG

MI EEG

ErrP EEG

ECoG

EMG

# Variety of Datasets

- Publicly available datasets
  - QMNIST: https://github.com/facebookresearch/qmnist
  - Stress: https://physionet.org/content/noneeg/1.0.0/
  - RSVP: http://hdl.handle.net/2047/D20294523
  - MI: https://physionet.org/physiobank/database/eegmmidb/
  - ErrP: https://www.kaggle.com/c/inria-bci-challenge
  - Faces: https://exhibits.stanford.edu/data/catalog/zk881ps0522
  - Ninapro: https://zenodo.org/record/1000116#.XuIppS2z3OR

| Datasets | Modality | Dimension | Nuisance ($|S|$) | Labels ($|Y|$) | Samples |
|---|---|---|---|---|---|
| QMNIST | Image | $28 \times 28$ | 539 | 10 | 60,000 |
| Stress | Temperature etc. | 7 | 20 | 4 | 24,000 |
| RSVP | EEG | $16 \times 128$ | 10 | 4 | 41,400 |
| MI | EEG | $64 \times 480$ | 106 | 4 | 9,540 |
| ErrP | EEG | $56 \times 250$ | 27 | 2 | 9,180 |
| Faces Basic | ECoG | $31 \times 400$ | 14 | 2 | 4,100 |
| Faces Noisy | ECoG | $39 \times 400$ | 7 | 2 | 2,100 |
| Ninapro | EMG | 16 | 10 | 12 | 890,446 |

# AutoBayes Benefit: Explore Different Models for Different Problems

- *No Free Lunch Theorems*: There is no one model that performs best for every dataset

Ensemble Methods in Machine Learning

# Ensemble Learning

- Every single model may be weak
- **Combining multiple weak models** may beat one strong model





*Tiny* weak fishes

*Gigantic* strong fish

# Ensemble Aggregation

- AutoML/AutoBayes explores many models

- Wasting by throwing away all weaker models?

- Employ **ensemble method** across explored models
  - Logistic regression (LR)
  - MLP
  - Transformer (multi-head attention)
  - ...



Final Class Score

Meta Learner
(Ensemble Fusion)

Ensemble Class Scores

Stacking Base Learners

Graphical Models

# Ensemble Leaning Gain in AutoBayes (QMNIST)

- Significant gain by ensemble learning; **1.3% gain**, state-of-the-art accuracy

# Ensemble Leaning Gain in AutoBayes (heartrate, EEG)

- Significant gain by ensemble learning; Up-to **37% gain**

# Ensemble Leaning Gain in AutoBayes (ECoG, EMG)

- Significant gain by ensemble learning; Up-to **12% gain**

Ensemble



Faces
Basic
68.83% → 78.36%

Faces
Noisy
77.50% → 89.71%

Ninapro
31.33% → 43.20%

# AutoBayes Ensemble Gain

# Subject Variation Robustness (Stress Dataset)

# Summary

- We introduced a new concept called **AutoBayes** for macro DNN architecture exploration
  - Different **Bayesian graphs** are explored systematically
  - **Bayes Ball** algorithm justifies pruning independent edges
  - Encoder, decoder, estimator, classifier, and adversary network blocks are rationally linked

- We also discussed **transfer learning, adversarial learning, ensemble learning**
  - Multiple architectures explored in AutoML are not wasted for final classification (as **base learners**)
  - Different **meta learners** (LR, MLP, Transformer) are evaluated to aggregate multiple models

- Demonstrated the benefit for various public physiological datasets
  - Various different modalities (**image, heartrate, EEG, ECoG, EMG**) and dimensionalities are considered

- Questions?
  - Contact us: koike@merl.com, yewang@merl.com
  - More details in arXiv: https://arxiv.org/abs/2007.01255

**DeepAI**

- https://deepai.org/publication/autobayes-automated-inference-via-bayesian-graph-exploration-for-nuisance-robust-biosignal-analysis

- When our arXiv was uploaded on July 2, it became *top trending paper*
  - Obtained **71 "likes" in 4 days**
  - It was highlighted in **"This Week in A.I." Newsletter** on July 11

# Contributions

- AutoBayes explores potential graphical models inherent to the data, rather than exploring hyperparameters of DNN blocks.

- AutoBayes offers a solid reason of how to connect multiple DNN blocks to impose conditioning and adversary censoring for the task classifier, feature encoder, decoder, nuisance indicator and adversary networks, based on an explored Bayesian graph.

- It provides a systematic automation framework to explore different inference models through the use of the Bayes-Ball algorithm and ordered factorization.

- The framework is also extensible to multiple latent representations and multiple nuisance factors.

- Besides fully-supervised training, AutoBayes can automatically build some relevant graphical models suited for semi-supervised learning.

- Ensemble learning is introduced to improve performance while AutoBayes model exploration

# AutoBayes as AutoML: Macro to Micro Exploration

- In principle, AutoBayes can be applied to arbitrary number of nodes

- Splitting X, Y, Z, S macro-nodes into scalar-valued micro-nodes, AutoBayes operates like AutoML architecture search but with more theoretical justification



(d) Model D    (k) Model K

Micro Nodes
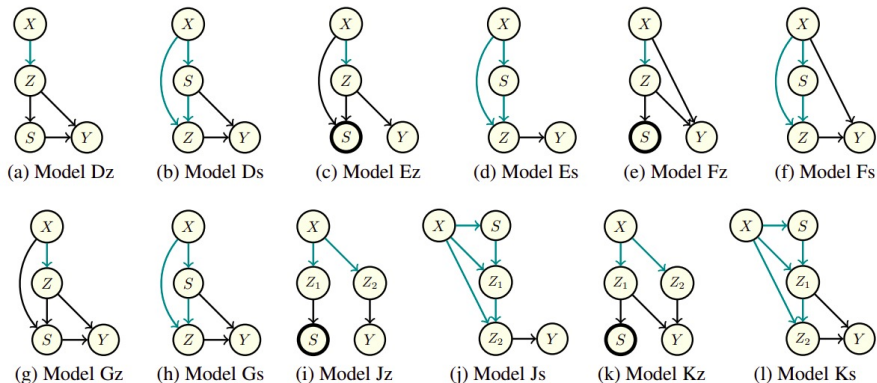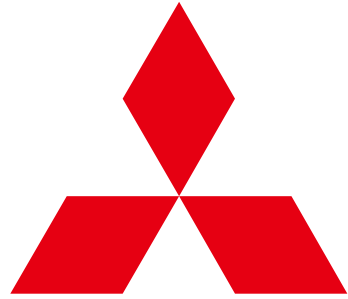
B

Define directed super-structure

A

AutoBayes

Use
Struc

Group using
Constraint Graph

C

Macro Nodes

# Methodology

$$p(y,s,z,x) = \begin{cases} p(y)p(s|\cancel{y})p(z|\cancel{s},y)p(x|\cancel{z},\cancel{s},y), & \text{Model-A} \\ p(y)p(s|\cancel{y})p(z|\cancel{s},y)p(x|z,\cancel{s},\cancel{y}), & \text{Model-B} \\ p(y)p(s|\cancel{y})p(z|\cancel{s},y)p(x|\cancel{z},s,y), & \text{Model-C} \\ p(y)p(s|\cancel{y})p(z|s,y)p(x|z,\cancel{s},\cancel{y}), & \text{Model-D} \\ p(y)p(s|\cancel{y})p(z|\cancel{s},y)p(x|z,s,\cancel{y}), & \text{Model-E} \\ p(y)p(s|\cancel{y})p(z|s,y)p(x|z,\cancel{s},y), & \text{Model-F} \\ p(y)p(s|\cancel{y})p(z|s,y)p(x|z,s,\cancel{y}), & \text{Model-G} \\ p(y)p(s|\cancel{y})p(z|s,y)p(x|z,\cancel{s},y), & \text{Model-H} \\ p(y)p(s|\cancel{y})p(z|s,y)p(x|z,s,y), & \text{Model-I} \\ p(y)p(s|\cancel{y})p(z_1|s,y)p(z_2|\cancel{z_1},\cancel{s},y)p(x|z_2,z_1,\cancel{s},y), & \text{Model-J} \\ p(y)p(s|\cancel{y})p(z_1|s,y)p(z_2|z_1,\cancel{s},y)p(x|z_2,z_1,\cancel{s},y), & \text{Model-K} \end{cases}$$



(a) Model Dz  (b) Model Ds  (c) Model Ez  (d) Model Es  (e) Model Fz  (f) Model Fs

(g) Model Gz  (h) Model Gs  (i) Model Jz  (j) Model Js  (k) Model Kz  (l) Model Ks

- Slash-cancelled factors from the full-chain case explicitly indicate independence.

- Conditional independence enables pruning links in the inference factor graphs.

$$p(y,z_1,z_2,s|x) = \begin{cases} p(z_1,z_2|x)p(y,s|z_1,z_2,\cancel{x}), & \text{z/ys} \\ p(z_1,z_2|x)p(s|z_1,\cancel{z_2},\cancel{x})p(y|\cancel{s},\cancel{z_1},z_2,\cancel{x}), & \text{Z-Inference} \\ p(z_1|x)p(z_2|z_1,x)p(s|z_1,\cancel{z_2},\cancel{x})p(y|\cancel{s},\cancel{z_1},z_2,\cancel{x}), & \text{z2/z1/s/y} \\ p(z_2|x)p(z_1|z_2,x)p(s|z_1,\cancel{z_2},\cancel{x})p(y|\cancel{s},\cancel{z_1},z_2,\cancel{x}), & \text{z1/z2/s/y} \\ p(z_1|x)p(s|z_1,x)p(z_2|s,z_1,x)p(y|\cancel{s},\cancel{z_1},z_2,\cancel{x}), & \text{z2/s/z1/y} \\ p(s|x)p(z_1|s,x)p(z_2|s,z_1,x)p(y|\cancel{s},\cancel{z_1},z_2,\cancel{x}), & \text{s/z2/z1/y} \\ p(z_1|x)p(s|z_1,\cancel{x})p(z_2|\cancel{s},z_1,x)p(y|\cancel{s},z_2,\cancel{z_1},\cancel{x}), & \text{z1/s/z2/y} \\ p(s|x)p(z_1|s,x)p(z_2|\cancel{s},z_1,x)p(y|\cancel{s},z_2,\cancel{z_1},\cancel{x}), & \text{s/z1/z2/y} \\ p(s|x)p(z_1,z_2|s,x)p(y|\cancel{s},z_2,\cancel{z_1},\cancel{x}), & \text{S-Inference} \\ \dots \end{cases}$$

(30)

Z-first and S-first inference graph models relevant for generative models D–G, J, and K