# Challenges in Machine Learning Research

## Masashi Sugiyama

Director, RIKEN Center for
Advanced Intelligence Project (AIP)

Professor, the University of Tokyo

RIKEN

AIP

東京大学
THE UNIVERSITY OF TOKYO
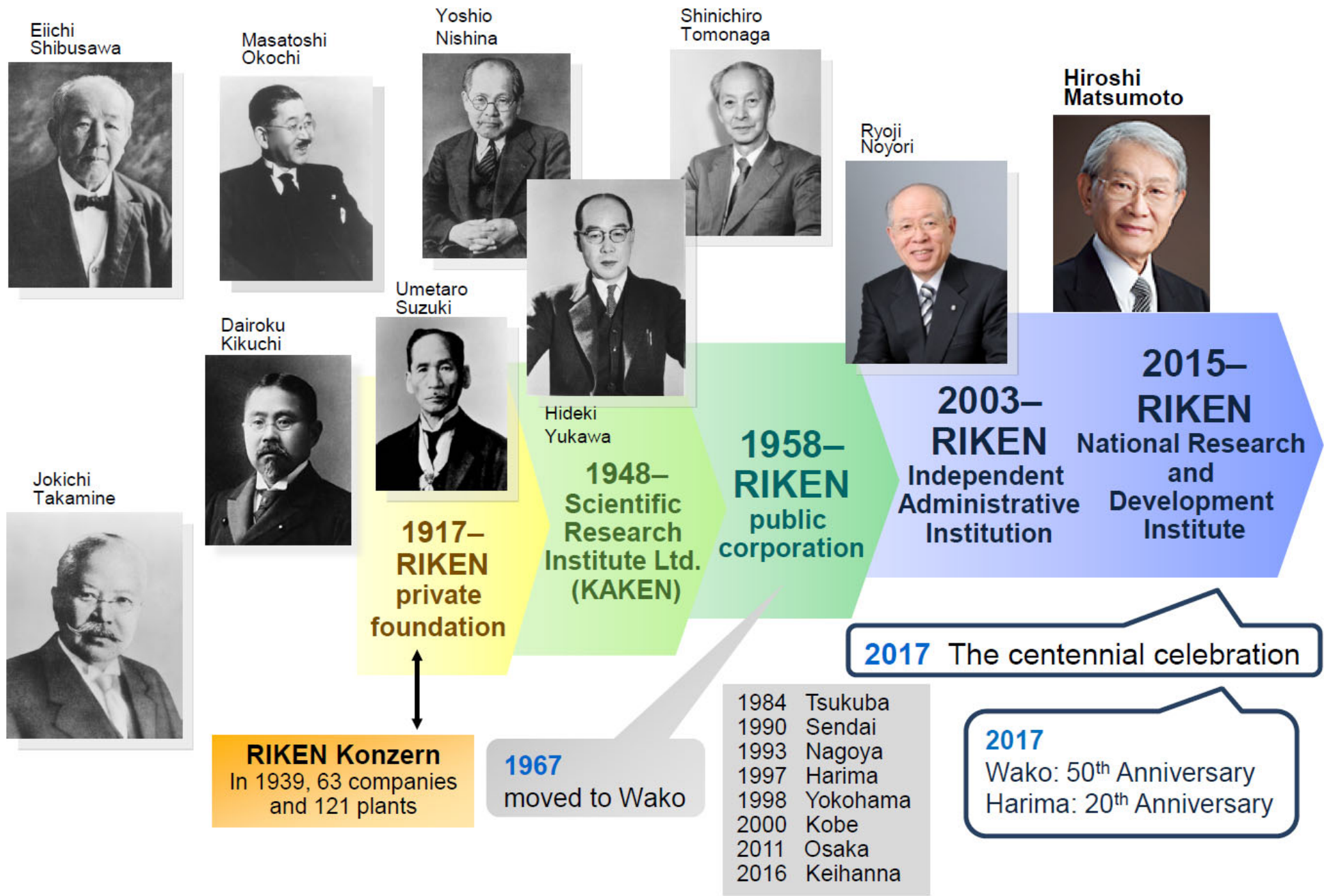
# What is "RIKEN"?

■ Name in Japanese:　理化学研究所

- Pronounced as:　　　rikagaku　kenkyusho
- Meaning:　Physics and Chemistry　Research Institute
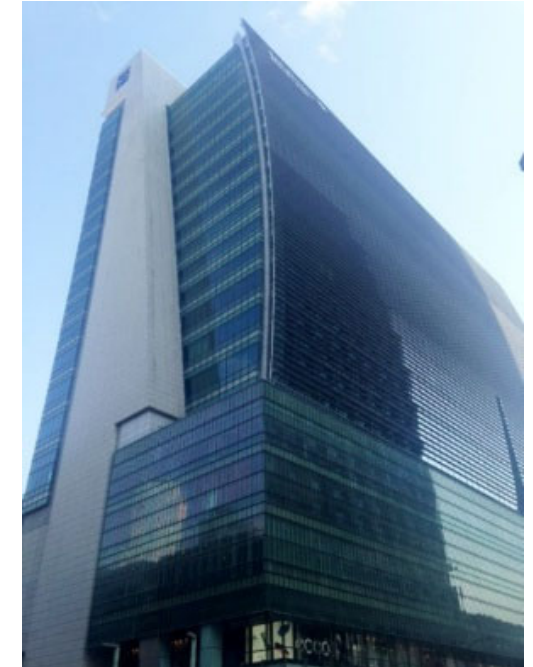
■ Acronym in Japanese: 理研 (RIKEN)

# Brief History

Eiichi Shibusawa

Masatoshi Okochi

Yoshio Nishina

Shinichiro Tomonaga

Hiroshi Matsumoto

Ryoji Noyori

Dairoku Kikuchi

Umetaro Suzuki

Hideki Yukawa

Jokichi Takamine

**1917– RIKEN private foundation**

**1948– Scientific Research Institute Ltd. (KAKEN)**

**1958– RIKEN public corporation**

**2003– RIKEN Independent Administrative Institution**

**2015– RIKEN National Research and Development Institute**

**RIKEN Konzern**
In 1939, 63 companies and 121 plants

**1967** moved to Wako

| | |
|---|---|
| 1984 | Tsukuba |
| 1990 | Sendai |
| 1993 | Nagoya |
| 1997 | Harima |
| 1998 | Yokohama |
| 2000 | Kobe |
| 2011 | Osaka |
| 2016 | Keihanna |

**2017** The centennial celebration

**2017**
Wako: 50th Anniversary
Harima: 20th Anniversary

# Research

RIKEN Cluster for Science, Technology and Innovation Hub/ Preventive Medicine and Diagnosis Innovation Program

Director Dr. Hidetoshi Kotera

Headquarters

President Dr. Hiroshi Matsumoto

Executive Directors

Dr. Shigeo Koyasu

Mr. Shigeharu Kato

Dr. Hidetoshi

Dr. Michihiko

Dr. Yuko Harayama

Nishina Center for

Spring-8 Center

Director Dr. Hiroshi Okaz

RIKEN Information R&D and Strategy Headquarters

RIKEN Cluster for Pioneering Research

BioResource Research Center

Director Dr. Yoshinori

Director Dr. Hiroyoshi

Director Dr. Tetsuya Ishikawa

Dr. Michihiko Mino

Dr. Shigeo Koyasu

Medicine

Physics

Physics

Biology

Physics

Neuro-science

Physics

Chemistry

Biology

Physics

Physics

Medicine

Informatics

Informatics

<Overseas>
RAL-RIKEN (UK)
BNL-RIKEN (USA)
Beijing Representative Office (China)
Singapore Representative Office (Singapore)
European Representative Office (Belgium)

Center for Computatio Science

Sendai
Tsukuba
Harima
Keihanna
Nagoya
Wako
Yokohama
Tokyo
Osaka

Director Katsumi Midorikawa

Director Dr. Ryuichiro Kageyama

Director Dr. Yasunobu Nakamura

Center f Integra Medical

Center for Advanced I oject

Director Dr. Eisuke Nishida

Director Dr. Kazuhiko Yamamoto

Director Dr. Kazuki Saito

Director Dr. Tetsuo Hatsuda

# What is RIKEN-AIP?

■ RIKEN founded Center for Advanced Intelligence Project (AIP) in 2016, under Ministry of Education, Culture, Sports, Science and Technology (MEXT).

MEXT
MINISTRY OF EDUCATION,
CULTURE, SPORTS,
SCIENCE AND TECHNOLOGY-JAPAN

Main office located in the heart of Tokyo



Distributed office across Japan



Sendai
Tsukuba
Kyoto  Shiga
Osaka
Fukuoka
Nara
Tokyo

In-house GPU servers



Open discussion space

# AIP's 5 Missions

■ Develop next-generation AI technology:

- machine learning & optimization theory, etc.

■ Accelerate scientific research:

- cancer, material, genomics, etc.

■ Solve socially critical problems:

- natural disaster, elderly healthcare, etc.

■ Study of ethical, legal and social issues of AI:

- ethical guidelines, personal data, etc.

■ Human resource development:

- researchers, engineers, etc.

# Statistics

As of Apr. 1, 2021

■ **Diverse research staffs**:

- 140 employed researchers
  (30% international, 20% female)
- 290 visiting researchers
- 60 domestic students
- 140 international interns (total)

■ **Extensive collaboration**:

- 3 industry collaborative centers
- 40+ industry projects
- 40+ international collaboration partners



**NEC**

**FUJITSU**

**TOSHIBA**

# Contents

1. Introduction
2. Research at RIKEN-AIP
   A) ML Application
   B) ML Society
   C) ML Theory
3. Research at IIL Team
4. Future Challenges

# AIP's Research Challenges

- **Machine learning (ML)** is the core of current AI:
  - Let a computer learn like humans.
  - Successful in speech, image, language, ads,…
- However, current ML is:
  - data-hangry (requiring big labeled data for training),
  - black-box (less interpretable).
- Our challenges:
  - Develop new ML theory to overcome the limitations.
  - Explore new ML application beyond current ML.
  - Design new ML society with appropriate ethical discipline and data-circulation systems.

# Contents

1. Introduction
2. Research at RIKEN-AIP
   A) ML Application
   B) ML Society
   C) ML Theory
3. Research at IIL Team
4. Future Challenges

# 1-1) Prostate Cancer Diagnosis

- **Prostate cancer** accounts for 10% of male cancers:
  - Automatic diagnosis is desired.
- Supervised classification needs **annotated** pathological images:
  - Increasing doctors' burden.
- Let's use **unsupervised** deep learning for feature extraction.

部位別がん罹患率
（全年齢）
［男性　2017年］
人口10万人対

| | |
|---|---|
| 口腔・咽頭, 25.0 | |
| 食道, 34.3 | |
| 胃, 144.9 | Stomach |
| 結腸, 88.2 | |
| 直腸, 53.0 | |
| 大腸, 141.1 | |
| 肝臓, 43.1 | |
| 胆のう・胆管, 19.6 | Colon |
| 膵臓, 34.4 | Lang |
| 喉頭, 7.9 | |
| 肺, 134.4 | |
| 皮膚, 19.7 | |
| 乳房, 1.1 | |
| 前立腺, 147.9 | |
| 膀胱, 28.0 | Prostate |
| 腎など, 32.5 | |
| 脳・中枢神経系, 5.0 | |
| 甲状腺, 7.5 | |
| 悪性リンパ腫, 30.0 | |
| 多発性骨髄腫, 6.8 | |
| 白血病, 13.0 | |

0  20  40  60  80  100  120  140  160  180

資料： 国立がん研究センターがん対策情報センター
Source: Center for Cancer Control and Information Services
National Cancer Center, Japan

# Unsupervised Deep Learning

Yamamoto et al. (Nature Communications 2019)
One of the top 50 most read Nature Communications articles in physics in 2019

■ We used 11+ billion unlabeled pathological image patches for feature extraction.



■ In addition to the standard Gleason score, novel features such as interstitium change were discovered.

■ Further applications in iPS cells, leukemia, and breast cancer.



AUC for 1-year recurrence prediction

# 1-2) Ghost Cytometry



Structured Illumination
for direct feature extraction based on optofluidics

Ota et al. (Science 2018)

Cell

FlowCytometry
1350x100:0.72 um/px

Machine Learning

PMT

FPGA

100-300μs (3000-10000cells/s)

Electrical Signal

Judge-ment in FPGA

Electrical Signal

Time

■ Classify normal/abnormal cells in the flow:

- However, deep learning inference is too slow.

■ Structured illumination allows direct feature extraction, resulting in real-time classification:

- Found a start-up for industrialization.
- Application in tumors and iPS cells.

- Nankai Trough is located south of Japan, expected to cause a big earthquake in the near future:
  - Risk assessment is indispensable.



Expected epicenter

Nankai Trough

©Google

https://www.fnn.jp/articles/-/22389

# Mathematical Model of Cycles

■There is a powerful mathematical model:

Equation of motion of ocean plate shear stress & land plate friction force

Land side plate

Ocean plate

Shear stress (pushing force)

$$\tau_i(t) = \sum_{j=1}^{N} K_{ij}\left(v_j^{pl}t - u_j(t)\right) - \frac{G}{2\beta}\frac{du_j(t)}{dt}$$

Friction force (stopping force)

○: Friction parameter

$$\tau_i(t) = \mu_i(t)\sigma_i^{eff}$$

$$\mu_i(t) = \mu_* + a_i \ln\left(\frac{V_i(t)}{V_*}\right) + b_i \ln\left(\frac{V_*\theta_i(t)}{L_i}\right)$$

■Tuning of friction parameters is the key.

● However, there are no enough supervised data.

Hachiya et al. (EGU2019)

■ **Alternately perform**

- **Simulation**: Generating artificial data by induction.
- **Learning**: Training a model with artificial data.

$(b_1^2 = 0.012, b_2^2 = 0.011)$

$(b_1^1 = 0.011, b_2^1 = 0.011)$

simulation

Learning

$V_2$

$V_1$

**Friction parameters**

**Earthquake cycle data**

■ **Prediction of earthquake cycles is highly improved.**

b1

Prediction error: 110 years

|True b - Predict b|

b2

Prediction error: 125 years

|True b - Predict b|

**Naïve method**

b1

Prediction error: 6.5 years

|True b - Predict b|

b2

Prediction error: 11 years
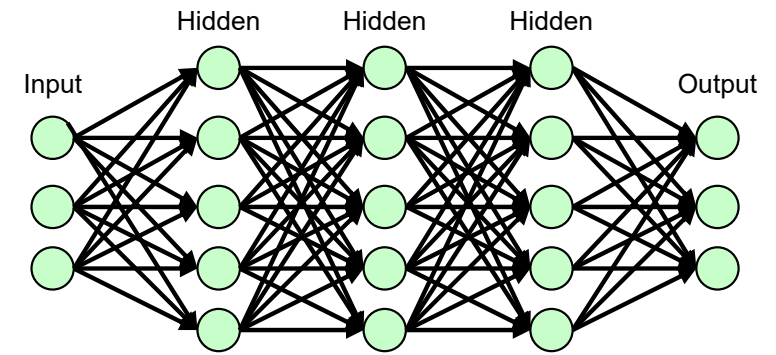
|True b - Predict b|

**Our method**

# Contents

1. Introduction
2. Research at RIKEN-AIP
   A) ML Application
   B) ML Society
   C) ML Theory
3. Research at IIL Team
4. Future Challenges

# 2-1) AI Ethical Guidelines

■ We have contributed to the discussions on privacy, fairness, security, etc.:

- Japanese Society for AI:
  - ■ Ethical Guidelines (2017).

- Ministry of Internal Affairs and Communications:
  - ■ AI R&D Guidelines, proposed to OECD (2017),
  - ■ AI Utilization Guidelines (2019).

- Cabinet Office:
  - ■ Social Principles of Human-centric AI, proposed to G20 (2019).

- IEEE:
  - ■ Ethically Aligned Design (2019).

# 2-2) Personal Life Repository

■How should personal information be managed?

- Company-based or government-based?

■We propose an individual-based system:

- Data subjects control data accessibility,
- Low-cost deployment.

■Proof-of-concept:

- Thousands of high schoolers share their learning records with the school management system.



personal device

personal app.

PLR

encrypte

PLR Cloud
(Google Drive,
Drpbox, etc.)

ジの活用

etc.)

encrypte

PLR

corporate app.

corporate computer

encrypted data

# Contents

1.  Introduction
2.  Research at RIKEN-AIP
    A)  ML Application
    B)  ML Society
    C)  ML Theory
3.  Research at IIL Team
4.  Future Challenges

# 3-1) Understanding Deep Learning [21]

🟧 Deep learning:

- Stacking many layers.

- Hard to optimization.

- Works excellently in practice.

🟧 We proved its superiority mathematically:

- Global optimization is possible.

   Suzuki & Akiyama (ICLR2021)

- Better prediction
   for high-dimensional data.

   Suzuki (NeurIPS2020), Nitanda & Suzuki (ICLR2021),
   Suzuki & Nitanda (NeurIPS2021)

- Universal approximator (INN).

   Teshima et al. (NeurIPS2020)

Hidden    Hidden    Hidden

Input                                Output

$$\widehat{L}(X_k) - \int \widehat{L}(x)\mathrm{d}\pi_\infty(x) \lesssim \exp\left(-\Lambda_\eta^* k\eta\right) + \frac{c_\beta}{\Lambda_0^*}\eta^{1/2-\kappa}$$

$$\mathrm{d}X_t = -\nabla\left(\widehat{L}(X_t) + \frac{\lambda}{2}\|X_t\|_{\mathcal{H}_K}^2\right)\mathrm{d}t + \sqrt{\frac{2}{\beta}}\mathrm{d}\xi_t$$

$$R_{\mathrm{lin}}(\mathcal{F}_\gamma) \gtrsim n^{-\frac{2\bar{\beta}+d}{2\bar{\beta}+2d}-\kappa'}$$

$$R_{\mathrm{lin}}(\mathcal{F}_\gamma) := \inf_{\widehat{f}:\mathrm{linear}} \sup_{f^\circ \in \mathcal{F}_\gamma} \mathbb{E}_{D_n}[\|\widehat{f} - f^\circ\|_{L_2(P_X)}^2]$$

$$\mathbb{E}_{D^n}\left[\mathbb{E}_{W_k}[\|f_{W_k} - f^\circ\|_{L_2(P_X)}^2 | D_n]\right] \lesssim n^{-\frac{\gamma}{\alpha_1 - 3\alpha_2 + 1}} + \Xi_k$$

# 3-2) Causal Inference

■ **Correlation** vs. **Causality**:

- The number of Nobel prize winners can be predicted by Chocolate consumption.

- But, eating more chocolate does not increase the number of Nobel prize winners.
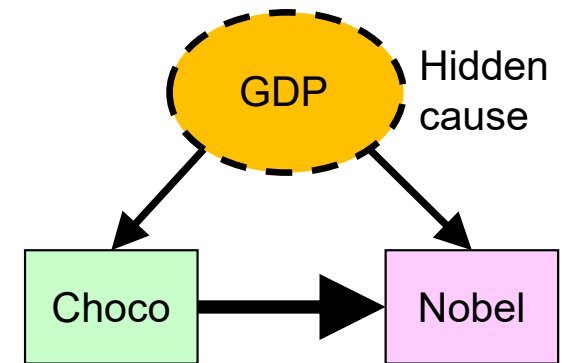
■ **Randomized controlled trial**:

- Split the subjects into two group.

- Treat only one and see what happens.

- Ethically problematic (e.g., vaccines)

Messerli (2012)

Corr: 0.79

Nobel Prize

Chocolate

Figure 1. Correlation between Countries' Laureates per 10 Million Population.

# Causal Inference in the Presence of Hidden Cause

■ In causal inference, how to handle <span style="color:red">hidden cause</span> is a big challenge!



GDP — Hidden cause

Choco → Nobel

■ <span style="color:red">We developed the first method to estimate the entire structure in the presence of hidden cause:</span>

- Speech separation technique is employed to separate hidden cause.

Maeda & Shimizu (AISTATS2020, UAI2021)

# Contents

1. Introduction
2. Research at RIKEN-AIP
3. Research at IIL Team
   A) Noisy-Label Learning
   B) Weakly Supervised Learning
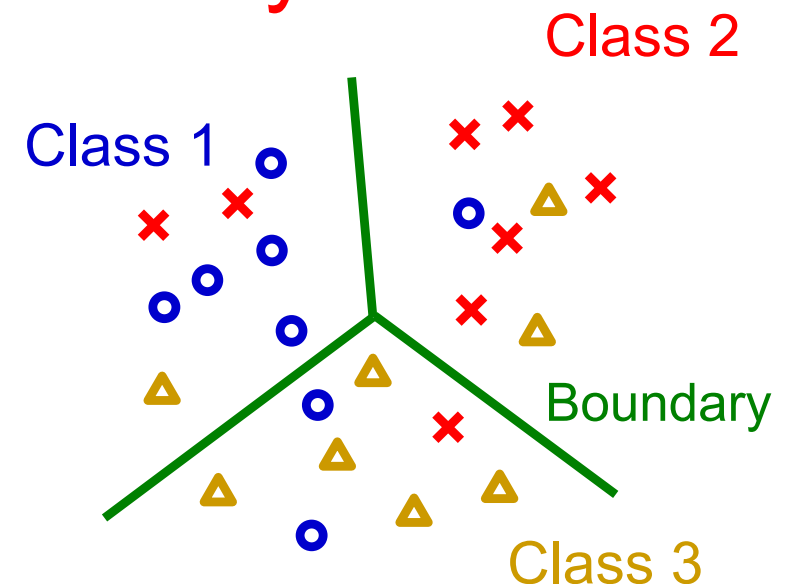   C) Transfer Learning
4. Future Challenges

# Imperfect Information Learning Team

- **Members**:
  - Gang Niu (Research Scientist): Learning theory
  - Shuo Chen (Postdoc): Metric learning
  - Jingfeng Zhang (Postdoc): Adversarial learning
  - Jiaqi Lyu (Postdoc): Weakly supervised learning
  - Many great Visiting Scientists,
    Junior Research Associates, Part-Timers,
    and Interns over the world!

- **Prof. Bo Han (HKBU) was an intern/postdoc:**
  - Now the most important Visiting Scientist in our team!

# Research at IIL Team

- **Goal**: Develop novel ML theories and algorithms that enable reliable learning from limited information.
  - **Label noise**: human error, sensor error.
  - **Insufficient information:** weak supervision.
  - **Data bias**: changing environments, privacy.
  - **Attack**: adversarial noise, distribution shift.

# Contents

# Supervised Classification

■ Supervised classification with clean labels:

Class 1

Class 2

Boundary

Class 3

Training error minimization is statistically consistent and work well in practice.

■ However, real-world labels are noisy possibly due to human error:

Training error minimization is no longer consistent and does not work well in practice.

Class 2

Class 1

Boundary

Class 3

# Generic Approaches

■ **Unsupervised outlier removal**:

- ● Substantially difficult

■ **Robust loss, regularization**:

- ● Not robust enough

■ We want to go beyond the limitations of existing approaches!

- ● Noise transition correction
- ● Clean sample selection

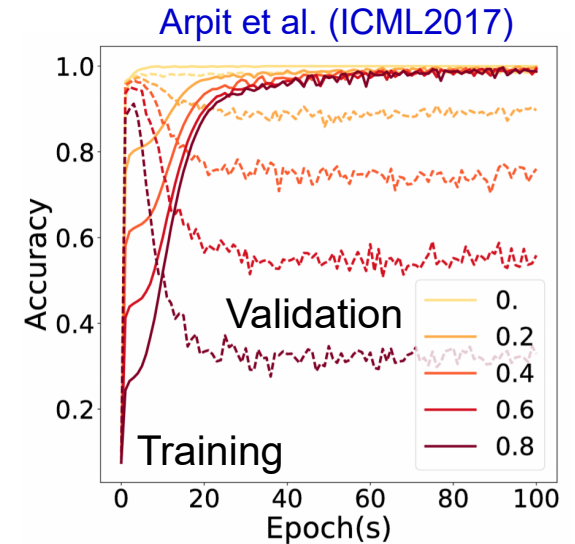# Noise Transition Correction

- **Noise transition matrix** $T$ :

  $$T^\top = \begin{array}{|c|c|c|} \hline 1 & 0.1 & 0.5 \\ \hline 0 & 0.8 & 0.5 \\ \hline 0 & 0.1 & 0 \\ \hline \end{array}$$

  - Clean-to-noisy flipping probability.

- Major approaches: Patrini et al. (CVPR2017)

  - **Loss correction** by $T^{-1}$ to eliminate noise.
  - **Classifier adjustment** by $T^\top$ to simulate noise.

- We want to estimate $T$ only from noisy data:

  - Use human cognition as a "mask" for $T$.
    Han et al. (NeurIPS2018)

  - Learn $T$ and a classifier dynamically.
    Xia et al. (NeurIPS2019)

  - Decompose $T$ into simpler components.
    Yao et al. (NeurIPS2020)

  - Regularize $T$ to be estimable.
    Zhang et al. (ICML2021), Li et al. (ICML2021)

  - Extension to input-dependent noise $T(x)$.
    Xia et al. (NeurIPS2020), Berthon et al. (ICML2021)

# Clean Sample Selection

■ Memorization of neural nets:

- Stochastic gradient descent fits clean data faster, but naïve early stopping does not work well.
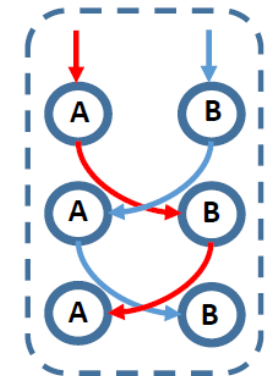
Arpit et al. (ICML2017)

■ "Co-teaching" with two neural nets:

- Teach small-loss data each other.
- Teach only disagreed data.
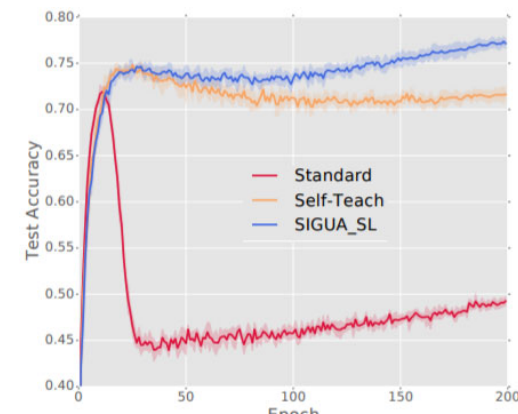- Gradient ascent for large-loss data.

Han et al. (NeurIPS2018)

Yu et al. (ICML2019)

Han et al. (ICML2020)

■ Very robust in experiments:

- Works well even if 50% of labels are randomly flipped.

# Contents

1. Introduction
2. Research at RIKEN-AIP
3. Research at IIL Team
   A) Noisy-Label Learning
   B) Weakly Supervised Learning
   C) Transfer Learning
4. Future Challenges

# Weakly Supervised Learning

■ Fully supervised data is expensive to collect.

■ **Weakly supervised data** can be collected easily:

- Ex.) Click prediction in online ads: It is easy to automatically collect
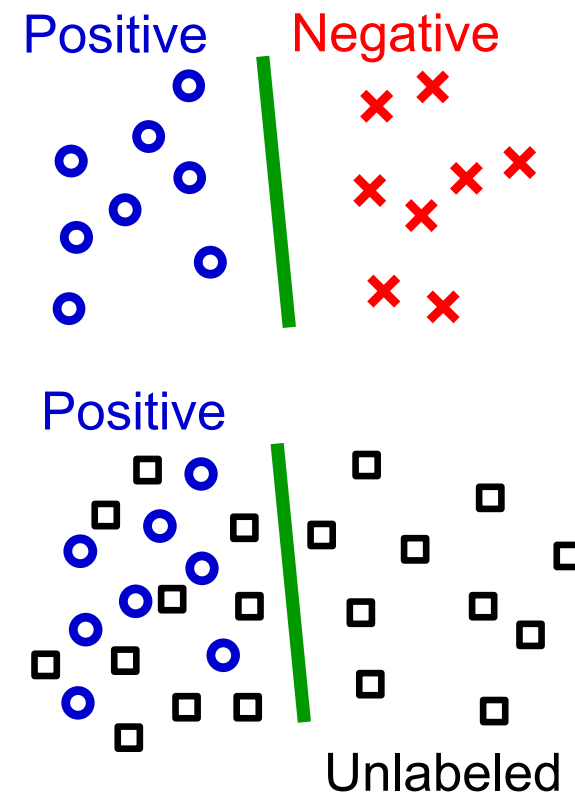  - Clicked ads (positive),
  - Unclicked ads (unlabeled).

Positive | Negative

Positive

Unlabeled

■ **Learning only from P and U data is possible!**

du Plessis et al. (NIPS2014, ICML2015, MLJ2017),
Niu et al. (NIPS2016), Kiryo et al. (NIPS2017), Hsieh et al. (ICML2019)

- Regard U data as noisy N data and correct the loss.
- Statistically consistent. $\mathcal{O}_p\left(1/\sqrt{n}\right)$
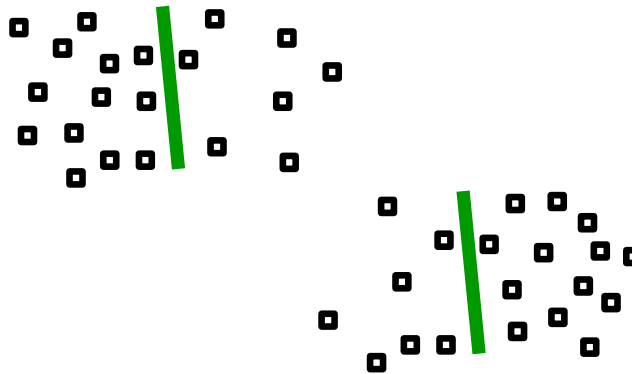
# Various Extensions

■ Learning from weakly supervised data is possible in many different forms!
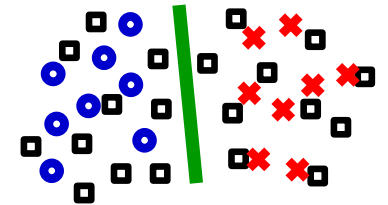
**Positive-Unlabeled**

du Plessis et al. (NIPS2014, ICML2015, MLJ2017)
Niu et al. (NIPS2016),, Kiryo et al. (NIPS2017)
Hsieh et al. (ICML2019)
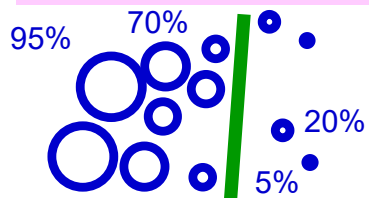
**Unlabeled-Unlabeled**

du Plessis et al.,(TAAI2013)
Lu et al. (ICLR2019, AISTATS2020)
Charoenphakdee et al. (ICML2019)
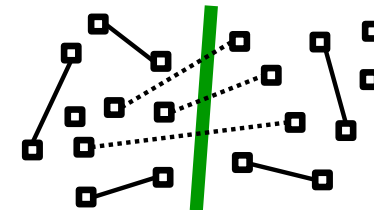Lei et al. (ICML2021)

**Semi-Supervised**

Sakai et al. (ICML2017, ML2018)

**Positive-confidence**

95%  70%
20%
5%

Ishida et al. (NeurIPS2018)
Shinoda et al. (IJCAI2021)

**Similar-Dissimilar**

Bao et al. (ICML2018)
Shimada et al. (NeCo2021)
Dan et al. (ECMLPKDD2021)
Cao et al. (ICML2021)
Feng et al. (ICML2021)

- All are loss-correction based and consistent. $\mathcal{O}_p\left(1/\sqrt{n}\right)$

- Any loss, classifier, and optimizer can be used.

# Multiclass Methods

■ Labeling patterns in multi-class problems is extremely painful.

■ Multi-class weak-labels:

- Complementary labels: Ishida et al. (NIPS2017, ICML2019) Chou et al. (ICML2020)
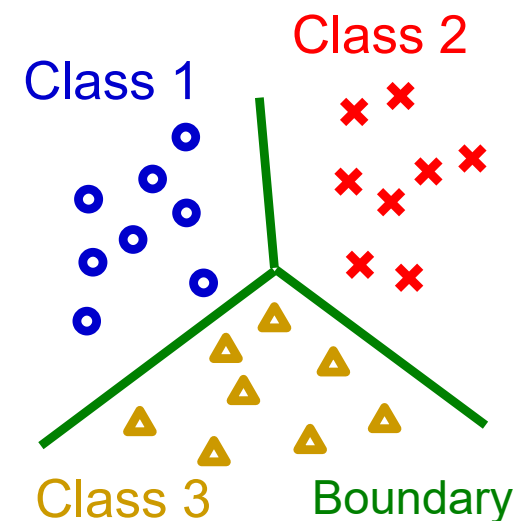  Specify a class that a pattern does not belong to ("not 1").

- Partial labels: Specify a subset of classes that contains the correct one ("1 or 2"). Feng et al. (ICML2020, NeurIPS2020) Lv et al. (ICML2020)

- Single-class confidence: Cao et al. (arXiv2021)
  One-class data with full confidence ("1 with 60%, 2 with 30%, and 3 with 10%")

■ Systematic loss correction is possible! $\mathcal{O}_p\left(1/\sqrt{n}\right)$

Class 1    Class 2

Class 3    Boundary

# Weakly Supervised Learning

**Supervised**

**Semi-supervised**

**Unsupervised**

P, N, U, S, D, Pconf, Nconf, Sconf, Dconf....
Comp, Partial, SCconf…
Different weak information can be systematically combined!

Labeling cost: High → Low

Low ← Classification accuracy → High

Sugiyama, Bao, Ishida, Lu, Sakai & Niu,
Machine Learning from Weak Supervision
MIT Press, 2022.

# Contents
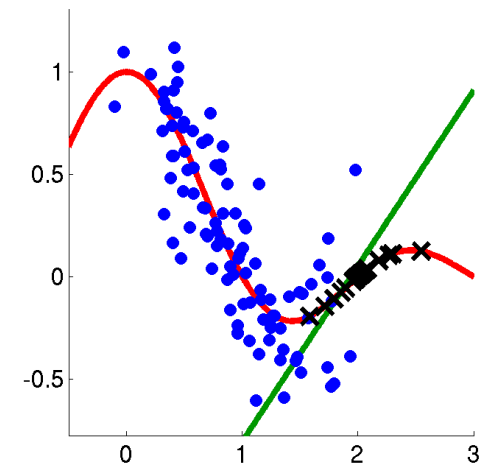
1. Introduction
2. Research at RIKEN-AIP
3. Research at IIL Team
   A) Noisy-Label Learning
   B) Weakly Supervised Learning
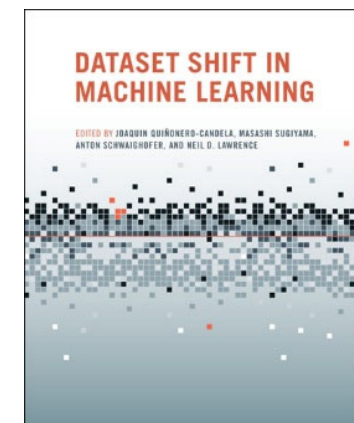   C) Transfer Learning
4. Future Challenges

# Transfer Learning

■ **Training and test data often have different distributions, due to**

- changing environments,
- sample selection bias (privacy).

■ **Transfer learning (domain adaptation):**

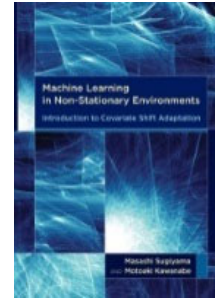- Train a test-domain predictor using training data from different domains.

NIPS Workshop 2006 - Whistler

NIPS Workshop on Learning when Test and Training Inputs Have Different Distributions, Whistler 2006

DATASET SHIFT IN MACHINE LEARNING

EDITED BY JOAQUIN QUIÑONERO-CANDELA, MASASHI SUGIYAMA, ANTON SCHWAIGHOFER, AND NEIL D. LAWRENCE

Quiñonero-Candela et al. (MIT Press 2009)

# Classical Approach for Transfer Learning

■ Two-step adaptation:

1. Importance weight estimation:

Sugiyama & Kawanabe
(MIT Press 2012)

$$\widehat{w} = \underset{w}{\arg\min} \, \widehat{\mathbb{E}}_{p_{\mathrm{tr}}(\boldsymbol{x},y)} \left[ D\left( w(\boldsymbol{x}, y), \frac{p_{\mathrm{te}}(\boldsymbol{x}, y)}{p_{\mathrm{tr}}(\boldsymbol{x}, y)} \right) \right]$$

2. Weighted predictor training:

$$\widehat{f} = \underset{f}{\arg\min} \, \widehat{\mathbb{E}}_{p_{\mathrm{tr}}(\boldsymbol{x},y)} [\widehat{w}(\boldsymbol{x}, y) \ell(f(\boldsymbol{x}), y)]$$

■ However, estimation error in Step 1 is not taken into account in Step 2.

● We want to integrate these two steps!

■ **Covariate shift:** Only input distributions change.

$$p_{\mathrm{tr}}(\boldsymbol{x}) \neq p_{\mathrm{te}}(\boldsymbol{x}) \qquad p_{\mathrm{tr}}(y|\boldsymbol{x}) = p_{\mathrm{te}}(y|\boldsymbol{x})$$

Shimodaira (JSPI2000)

■ Suppose we are given

- Labeled training data: $\{(\boldsymbol{x}_i^{\mathrm{tr}}, y_i^{\mathrm{tr}})\}_{i=1}^{n_{\mathrm{tr}}} \overset{\mathrm{i.i.d.}}{\sim} p_{\mathrm{tr}}(\boldsymbol{x}, y)$

- Unlabeled test data: $\{\boldsymbol{x}_i^{\mathrm{te}}\}_{i=1}^{n_{\mathrm{te}}} \overset{\mathrm{i.i.d.}}{\sim} p_{\mathrm{te}}(\boldsymbol{x})$

■ Minimize a risk upper bound jointly
w.r.t. weight $w$ and predictor $f$:

Zhang et al.
(ACML2020, SNCS2021)

$$J_{\ell_{\mathrm{tr}}}(f, w) \geq R_{\ell_{\mathrm{te}}}(f)^2$$

$$\widehat{f} = \operatorname*{argmin}_{f} \min_{w \geq 0} \widehat{J}_{\ell_{\mathrm{tr}}}(f, w)$$

$$R_\ell(f) = \mathbb{E}_{p_{\mathrm{te}}(\boldsymbol{x}, y)}[\ell(f(\boldsymbol{x}), y)]$$

$$\ell_{\mathrm{te}} \leq 1, \ell_{\mathrm{tr}} \geq \ell_{\mathrm{te}}$$

$\widehat{J}_\ell$ : Empirical approximation of $J_\ell$

- **Theoretical guarantee:**

$$R_{\ell_{\mathrm{te}}}(\widehat{f}) \leq \sqrt{2} \min_f R_{\ell_{\mathrm{te}}}(f) + \mathcal{O}_p(n_{\mathrm{tr}}^{-1/4} + n_{\mathrm{te}}^{-1/4})$$

# Dynamic Importance Weighting

- **General changing distributions:** $p_{\mathrm{tr}}(\boldsymbol{x}, y) \neq p_{\mathrm{te}}(\boldsymbol{x}, y)$

- **Suppose we are given**

  - Labeled training data: $\{(\boldsymbol{x}_i^{\mathrm{tr}}, y_i^{\mathrm{tr}})\}_{i=1}^{n_{\mathrm{tr}}} \overset{\mathrm{i.i.d.}}{\sim} p_{\mathrm{tr}}(\boldsymbol{x}, y)$

  - Labeled test data: $\{(\boldsymbol{x}_i^{\mathrm{te}}, y_i^{\mathrm{te}})\}_{i=1}^{n_{\mathrm{te}}} \overset{\mathrm{i.i.d.}}{\sim} p_{\mathrm{te}}(\boldsymbol{x}, y)$

- **For each mini-batch** $\{(\bar{\boldsymbol{x}}_i^{\mathrm{tr}}, \bar{y}_i^{\mathrm{tr}})\}_{i=1}^{\bar{n}_{\mathrm{tr}}}, \{(\bar{\boldsymbol{x}}_i^{\mathrm{te}}, \bar{y}_i^{\mathrm{te}})\}_{i=1}^{\bar{n}_{\mathrm{te}}}$,
  importance weights are estimated by
  matching **losses** by **kernel mean matching:**

  Fang et al.
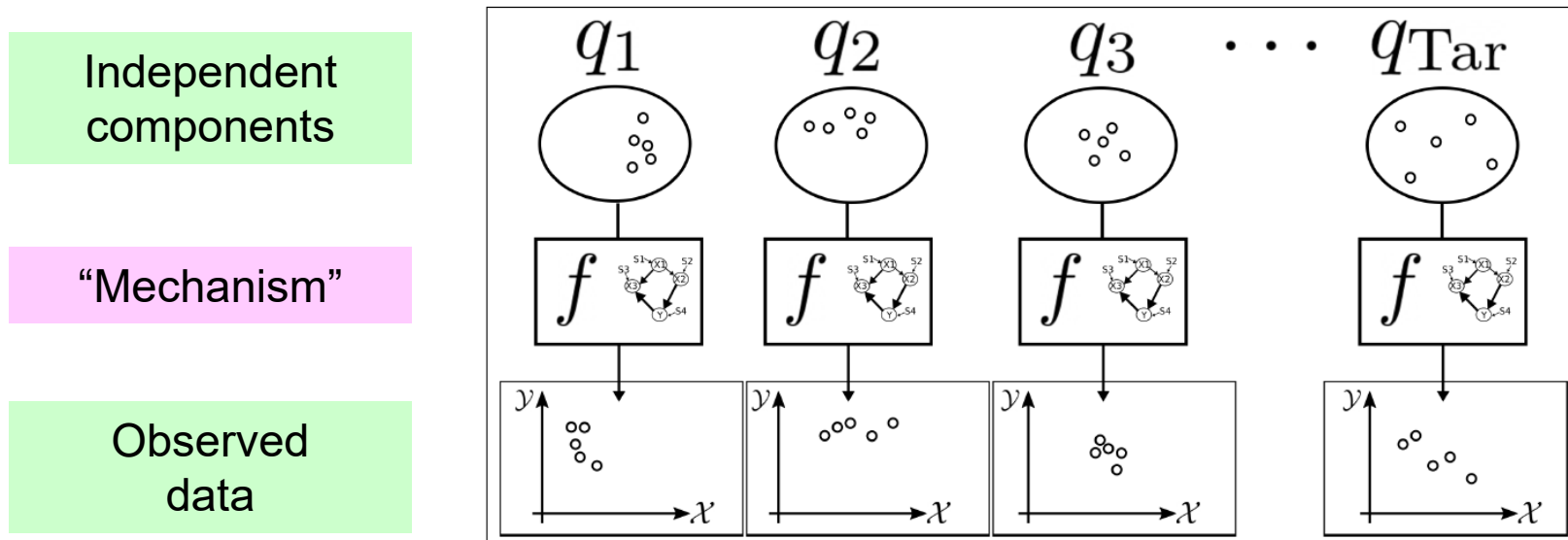  (NeurIPS2020)

  Huang et al. (NeurIPS2007)

$$\frac{1}{\bar{n}_{\mathrm{tr}}} \sum_{i=1}^{\bar{n}_{\mathrm{tr}}} r_i \ell(f(\bar{\boldsymbol{x}}_i^{\mathrm{tr}}), \bar{y}_i^{\mathrm{tr}}) \approx \frac{1}{\bar{n}_{\mathrm{te}}} \sum_{j=1}^{\bar{n}_{\mathrm{te}}} \ell(f(\bar{\boldsymbol{x}}_j^{\mathrm{te}}), \bar{y}_j^{\mathrm{te}})$$

- **Extremely simple, but highly powerful!**

# Mechanism Transfer

■ Is transfer learning possible when data distributions are seemingly very different?

■ Yes, if data generation mechanisms are shared:

- Use invertible neural networks (INNs) to invert the data generation mechanism.

  Teshima et al. (ICML2020)

- INNs are universal approximators.

  Teshima et al. (NeurIPS2020)

# Contents

1. Introduction
2. Research at RIKEN-AIP
3. Research at IIL Team
4. Future Challenges

# Challenges in Reliable ML

- **Reliability for expectable situations:**
    - Model the corruption process explicitly and correct the solution.
        - How to handle modeling error?

- **Reliability for unexpected situations:**
    - Consider worst-case robustness ("min-max").
        - How to make it less conservative?
    - Include human support ("rejection").
        - How to handle real-time applications?

- **Exploring somewhere in the middle would be practically more useful:**
    - Use partial knowledge of the corruption process.

# History of AI and Future

■ Classic AI:
- 1960s:
  symbolic, logical AI
- 1980s:
  Expert systems

■ Neuro-inspired AI:
- 1960s:
  1-layer perceptrons
- 1980s:
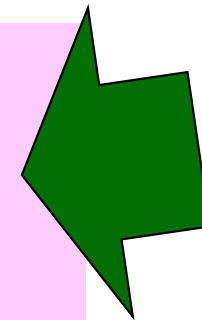  Multilayer perceptrons

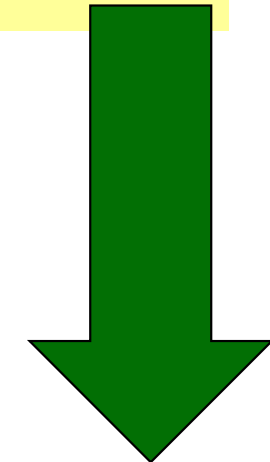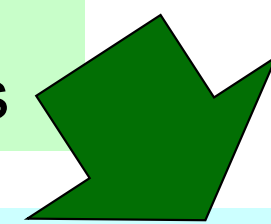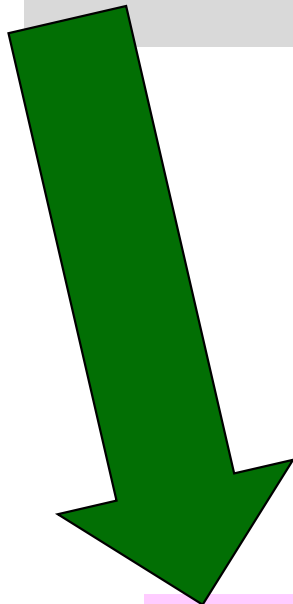■ Statistical machine learning:
- 2000s: Statistics, Bayes,
  convex optimization, kernels

■ Deep learning:
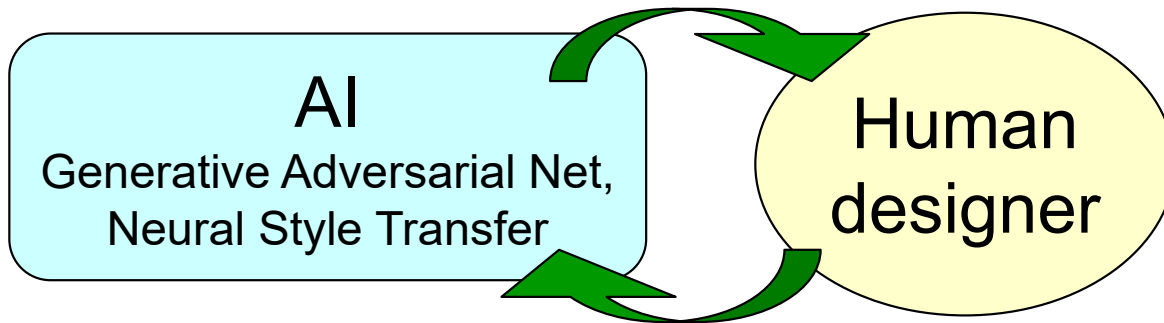- 2010s: Stochastic
  gradient, gigantic
  deep models

■ Next-generation AI:
- Integration of elements
- Human-like AI?

# Next-Generation AI

■ **Is Human-like AI ultimate?**

- Future AI needs not be autonomous.
- Future AI may learn together with humans.

AI
Generative Adversarial Net,
Neural Style Transfer

Human designer

■ **AI needs to be inclusive to human society:**

- Technology
  X
  Human creativity,
  culture, and ethics.

Fashion show at UTokyo in Mar. 2019
(with Prof. Aihara and Emarie)



https://www.fashion-press.net/collections/11006

# Thank You!
# 多謝！
# ありがとう！