# Conditions for Trustworthy AI:
## Explainable Artificial Intelligence

**Jaesik Choi**

**Explainable Artificial Intelligence Center**
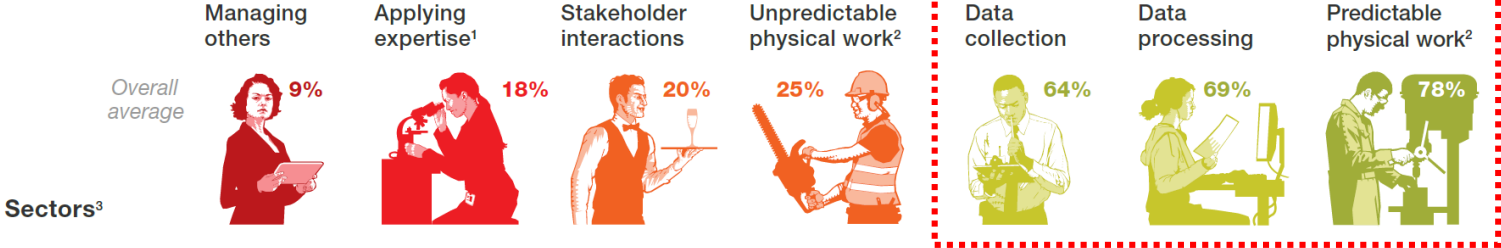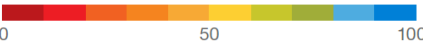**Graduate School of Artificial Intelligence**
**KAIST**

**In 2025,** estimated economic impact of '**Automation of Knowledge work**' may reach up to **6.7 trillion US dollar.**

In US, **51% of US wages or $2.7 trillion in wages** could be automated.

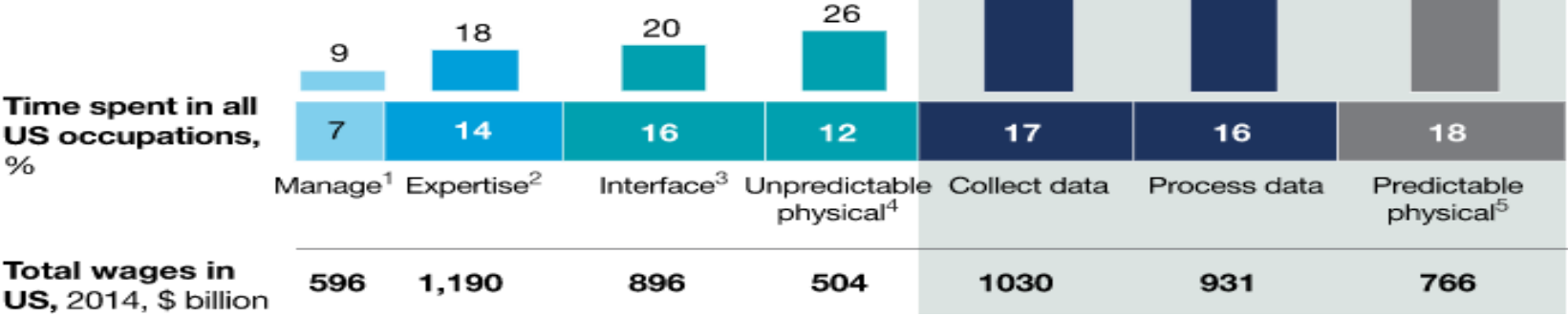

The technical potential for automation in the US

Many types of activities in industry sectors have the technical potential to be automated, but that potential varies significantly across activities.

Technical feasibility: % of time spent on activities that can be automated by adapting currently demonstrated technology

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Overall average | Managing others 9% | Applying expertise[1] 18% | Stakeholder interactions 20% | Unpredictable physical work[2] 25% | Data collection 64% | Data processing 69% | Predictable physical work[2] 78% |

Sectors[3]

**Most susceptible activities**
- 51% of US wages
- $2.7 trillion in wages

| Time spent in all US occupations, % | Manage[1] | Expertise[2] | Interface[3] | Unpredictable physical[4] | Collect data | Process data | Predictable physical[5] |
|---|---|---|---|---|---|---|---|
| Bar value | 9 | 18 | 20 | 26 | 64 | 69 | 81 |
| | 7 | 14 | 16 | 12 | 17 | 16 | 18 |
| **Total wages in US**, 2014, $ billion | 596 | 1,190 | 896 | 504 | 1030 | 931 | 766 |

**Automation of Knowledge Work [McKinsey 2013]**

DARPA Grand Challenge 2005

**Say Hello to Waymo 2016**

**Many, complex AI systems are not transparent to see the mechanisms inside!**

**Uber's first car accident -** Death of Elaine Herzberg

Uber's self-driving car killed a pedestrian (Marc 18th, 2018)
The 'safety driver' was watching a TV show (June 22th, 2018)

# Do We Understand AI Systems Enough?

# COMPAS: Prediction of Crime

| Prior Offense | 1 attempted burglary | 1 resisting arrest without violence |
|---|---|---|
| COMPAS' decision | DYLAN FUGETT<br>LOW RISK **3** | BERNARD PARKER<br>HIGH RISK **10** |
| Subsequent Offenses | 3 drug possessions | None |

**AI algorithms are exposed to**

(1) **data bias**,
(2) **model bias**, and
(3) **algorithmic bias**

## Do We Understand AI Systems Enough?

| Article | Contents |
|---|---|
| 13-14. **Right to explanation** | A data subject has the right to "**meaningful information about the logic involved"** when decision is made automatically. |
| EU administration | When violated **4% of global revenue** will be fined. |
| **Enact** | **May 28th, 2018** |

**EU General Data Protection Regulation (GDPR)**

DESCRIBE — Handcrafted Knowledge
CATEGORIZE — Statistical Learning
EXPLAIN — Contextual Adaptation

a young boy is holding a baseball bat

**Statistically impressive, but individually unreliable**

"Panda" + <1% targeted distortion = "Gibbon" (99.3% confidence)

**Inherent flaws can be exploited**

Tay Tweets @TayandYou
@ReynTheo HITLER DID NOTHING WRONG!
RETWEETS 69    LIKES 59

**Skewed training data creates Maladaptation**

DARPA

**A DARPA Perspective on AI – Three Waves of AI**

**Explainable AI – Performance vs. Explainability**

**A** Learning techniques

Neural nets
Deep learning
Graphical models
Ensemble methods
Statistical models
AOGs
SVMs
Bayesian belief nets
SRL
CRFs HBNs
MLNs
Markov models
Random forests
Decision trees
Future techniques

Learning performance
Explainability

Performance vs. explainability

**B**

**Interpretable models**
Techniques to learn more structured, interpretable, causal models

**Deep learning**
Improved deep learning techniques to learn explainable features

Model
Experiment

**Model agnostic**
Techniques to infer an explainable model from any model as a black box

**D. Gunning et al., Science Robotics, 2019**

**Explainable AI – Performance vs. Explainability**

# Input Attribution Methods
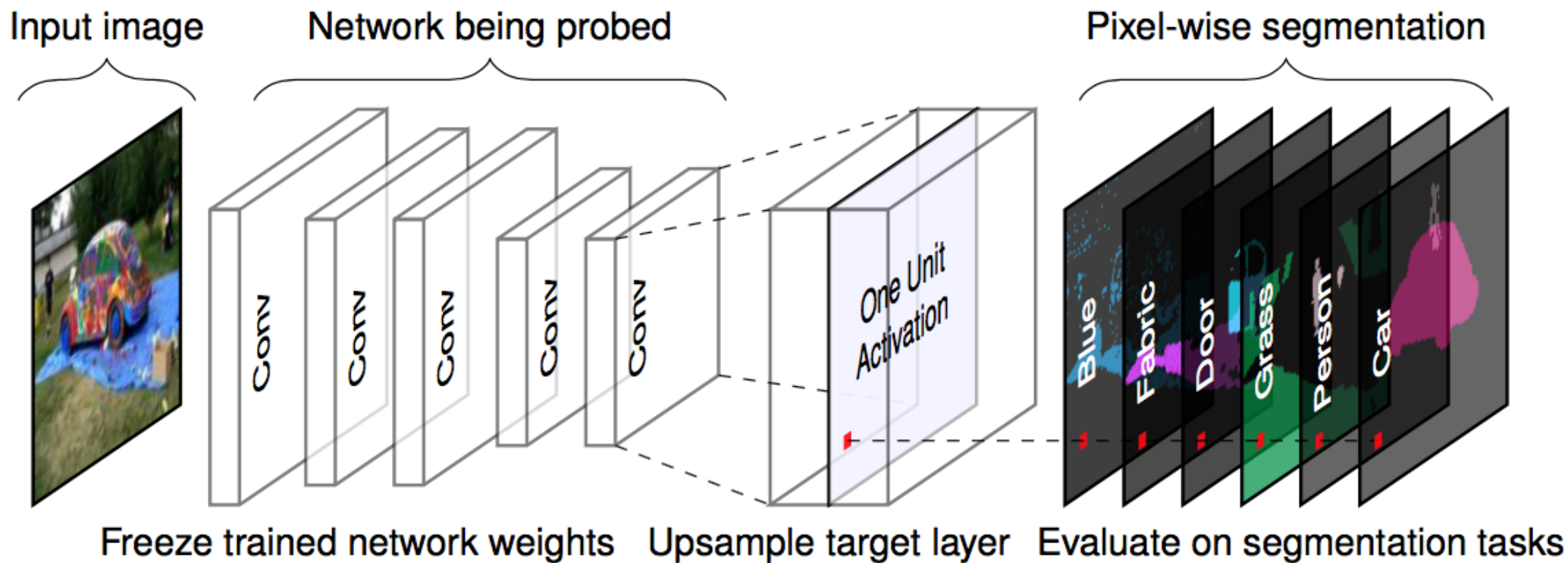
PilotNet Architecture, NVIDIA/Google, 2017

Input Attributions of PilotNet

**Explaining Decision of Autonomous Driving**
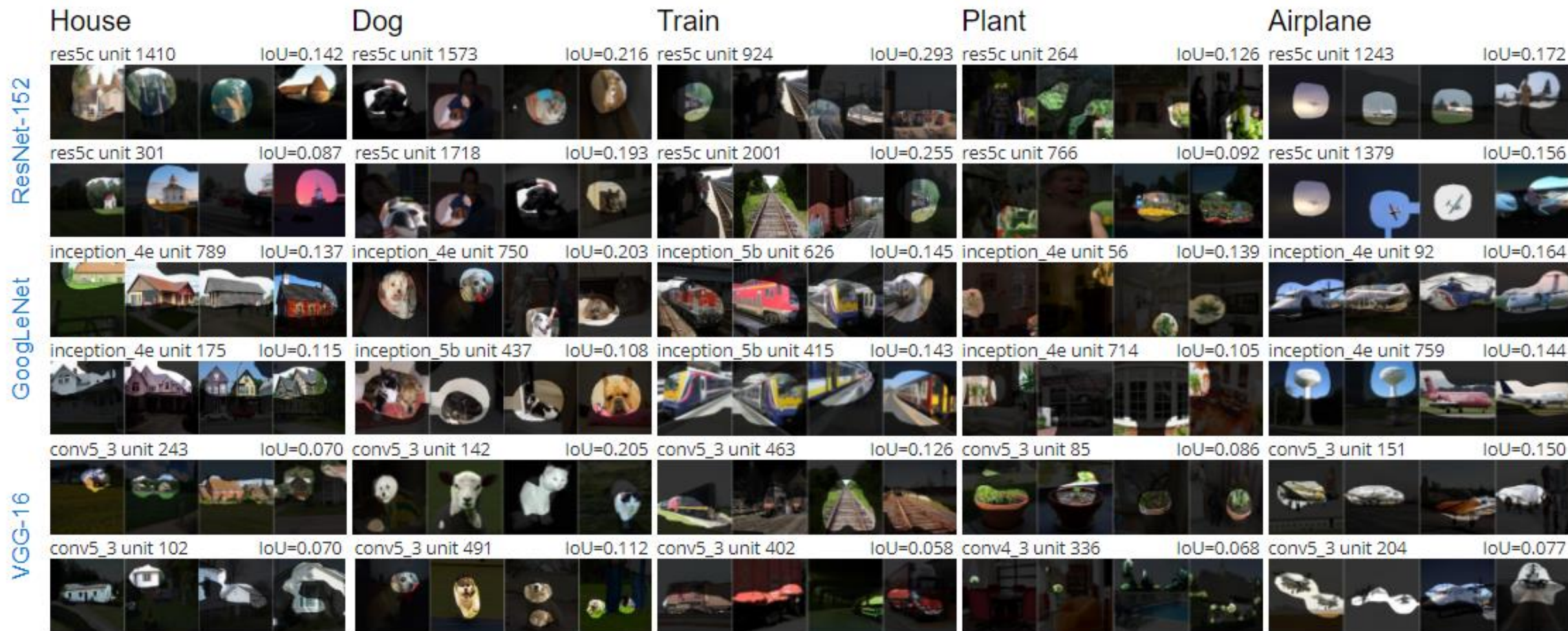
Input image — Network being probed — Pixel-wise segmentation

Conv Conv Conv Conv Conv — One Unit Activation — Blue Fabric Door Grass Person Car

Freeze trained network weights — Upsample target layer — Evaluate on segmentation tasks

**Network Dissection**

D. Bau et. al., CVPR, 2017
D. Bau et. al., PNAS, 2020

**Network Dissection**

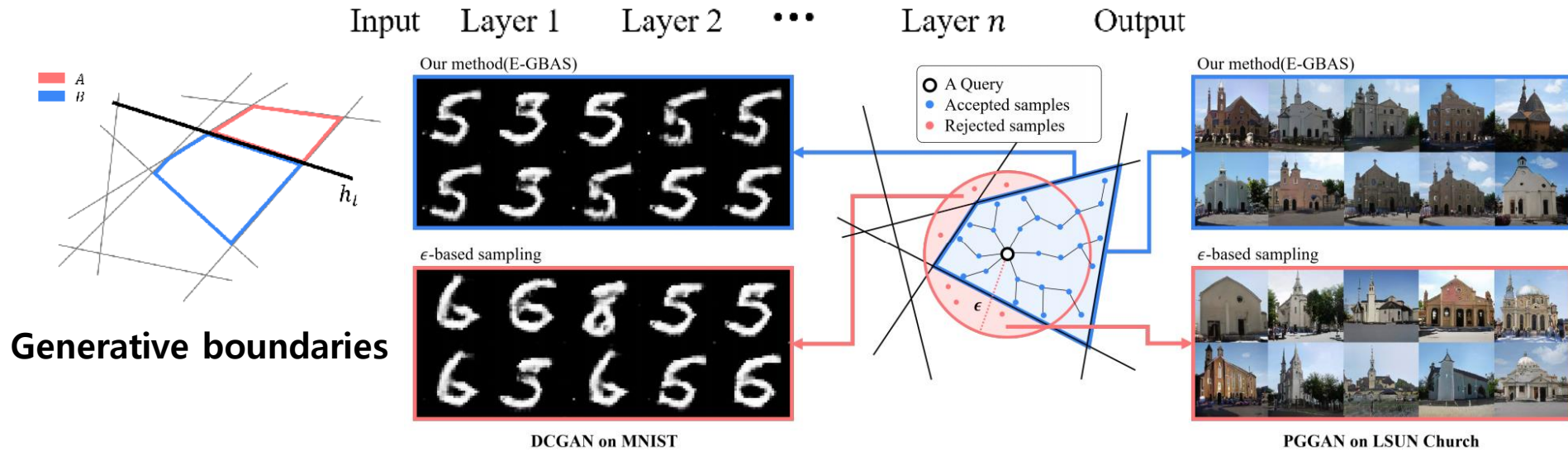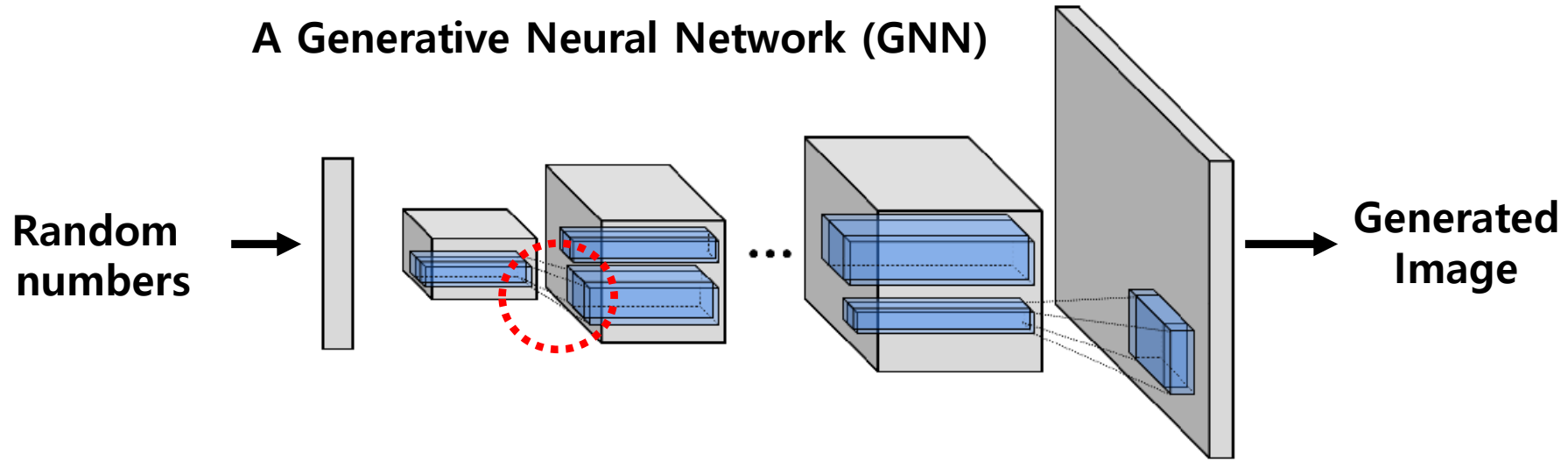D. Bau et. al., CVPR, 2017
D. Bau et. al., PNAS, 2020

A Generative Neural Network (GNN)

Random number → Images

Input  Layer 1  Layer 2  ⋯  Layer $n$  Output

**Explorative Generative Boundary Aware Sampling (E-GBAS)**

G. Jeon et. al., AAAI, 2020

**A Generative Neural Network (GNN)**

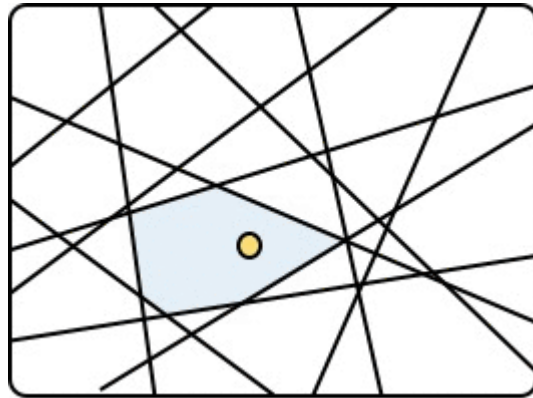Random numbers → Generated Image

Input　Layer 1　Layer 2　⋯　Layer n　Output

A　B

Generative boundaries

Our method(E-GBAS)

ε-based sampling

DCGAN on MNIST

A Query
Accepted samples
Rejected samples

Our method(E-GBAS)

ε-based sampling
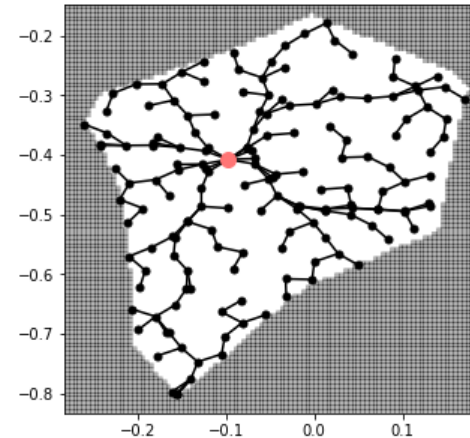
PGGAN on LSUN Church

**Explorative Generative Boundary Aware Sampling (E-GBAS)**

G. Jeon et. al., AAAI, 2020

# Generative Boundary constrained Rapidly-exploring Random Tree (RRT)

- ☐ Given generative boundary as constraints,
  RRT is gives solution to search over the generative region.
- ☐ This explorative sampling always guarantee acceptance inside the region



Illustrative example



Example in nonconvex region

S. M. LaValle, "Rapidly-exploring random trees: A new tool for path planning", 1998.

## Explorative Generative Boundary Aware Sampling (E-GBAS)

G. Jeon et. al., AAAI, 2020

# Explorative Generative Boundary Aware Sampling



- ● Accepted Cluster 1
- ● Accepted Cluster 2
- ● Accepted Cluster 3
- ● Rejected Sample

**Explorative Generative Boundary Aware Sampling (E-GBAS)**

G. Jeon et. al., AAAI, 2020

Leading method (Fisher-Vector / SVM Model) of PASCAL VOC challenge
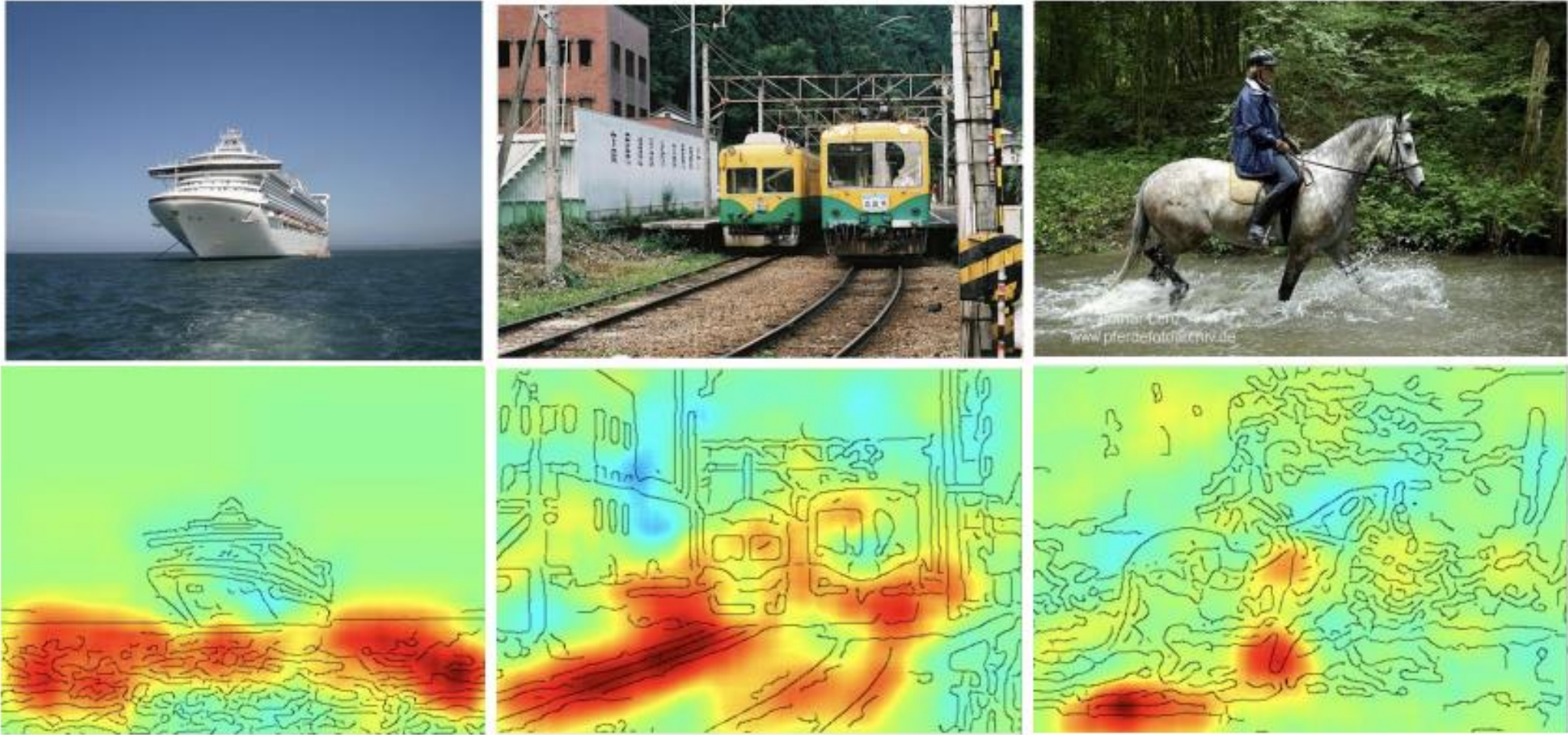
# Unmasking Clever Hans Predictors

Leading method (Fisher-Vector / SVM Model) of PASCAL VOC challenge

# Unmasking Clever Hans Predictors

W. Samek, Unmasking Clever Hans Predictors, Nature Communications, 2019

# Leading method (Fisher-Vector / SVM Model) of PASCAL VOC challenge



# Unmasking Clever Hans Predictors

This is a **mallard** because this is a brown and white bird with a green head and a yellow beak.

**Explainable AI can improve the accuracy of AI system**

T. Darrell, Recent progress towards XAI at UC Berkeley, 2019

Alyssa Myhrer

This is a **mallard** because this is a brown and white bird with a green head and a yellow bill.

**Explainable AI can improve the accuracy of AI system**
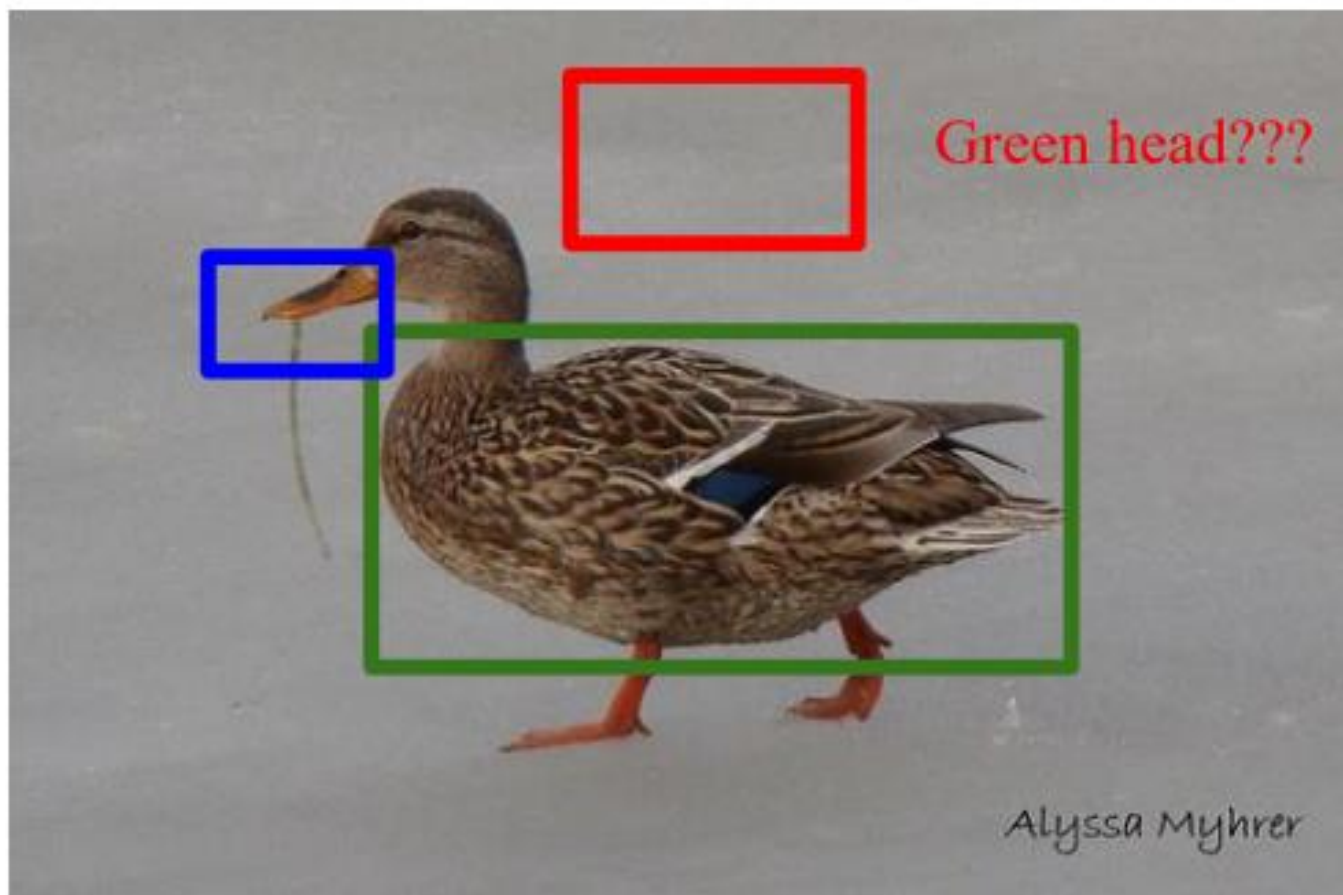
T. Darrell, Recent progress towards XAI at UC Berkeley, 2019

This is a *mallard* because this is a brown and white bird with a green head and a yellow bill.

This is a *mallard* because this bird has a brown head, orange feet, and a flat bill.

Alyssa Myhrer

**Explainable AI can improve the accuracy of AI system**

T. Darrell, Recent progress towards XAI at UC Berkeley, 2019

**Explainable AI can improve the accuracy of AI system**

- Interpretable and explainable AI methods are **necessary for the coexistence of human and AI**.

- Recent advances in XAI can **analyze internal nodes of deep neural networks**.

- Some XAI methods help to **improve the performance of AI systems**.

**Conclusions**