FUJITSU

shaping tomorrow with you

# Explainable AI that Can be Used for Judgment with Responsibility

FUJITSU LABORATORIES LTD.
25th April 2018

# About me

■ Name: Hajime Morita

■ Research interests:
  ■ Natural Language Processing
    • Summarization
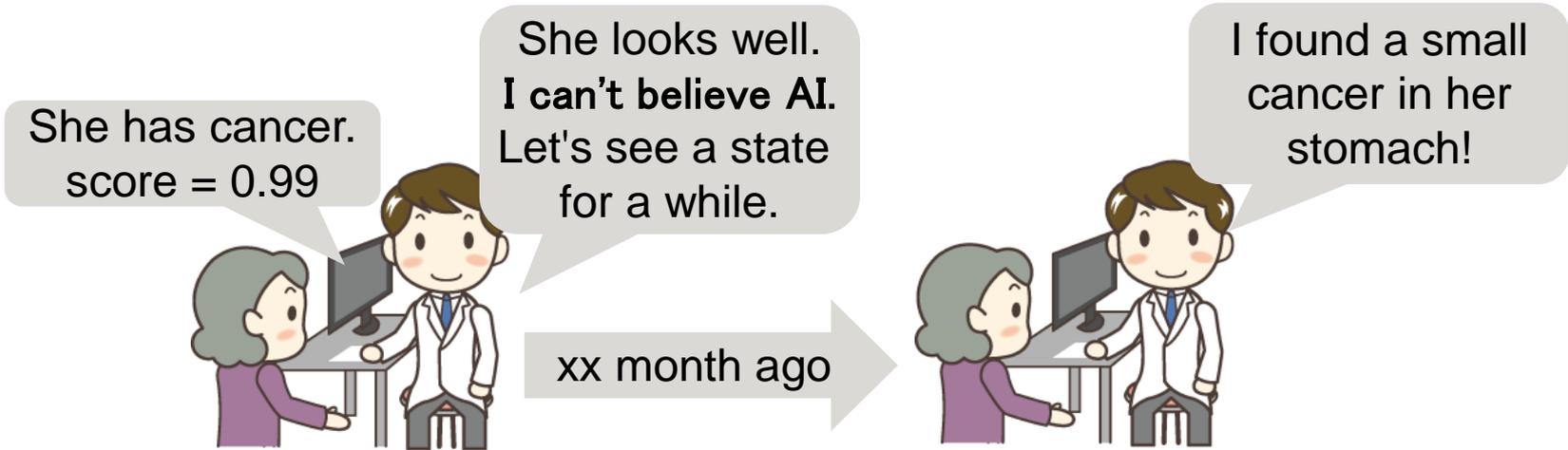    • Morphological Analysis
    • Information Extraction

■ Recent jobs:
  ■ I received my Ph.D. from Tokyo Institute of Technology (2015)
  ■ Researcher @ Kyoto University (2015 - 2017)
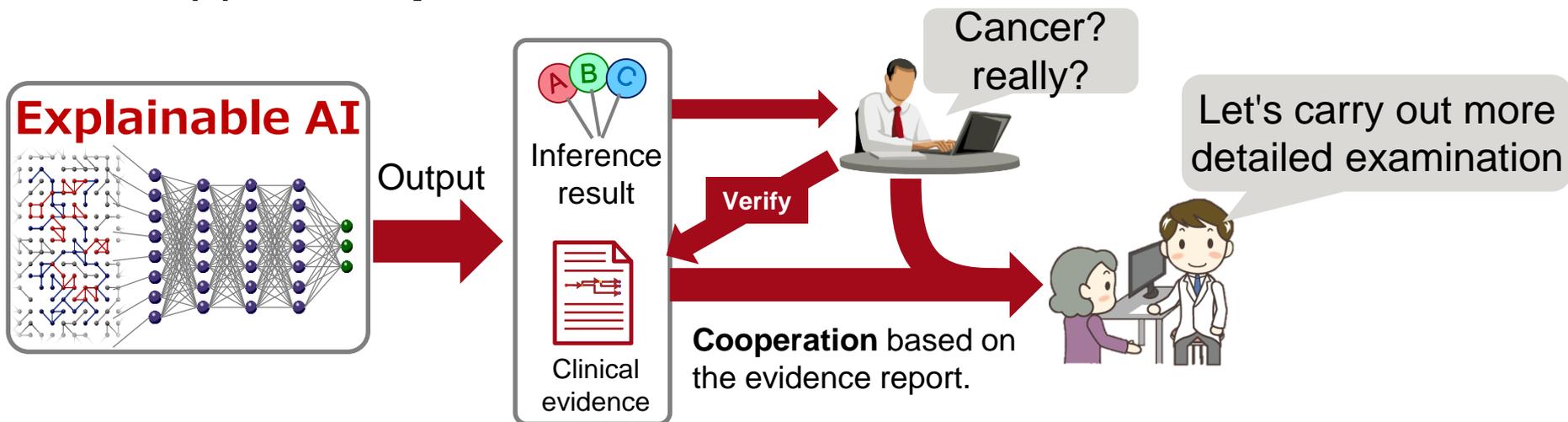  ■ Researcher @ Fujitu Laboratory (2017 -)

# Outline

- Introduction of explainable AI
- Knowledge Graph (short)
- Deep Tensor
- Explainable AI

# Why we need explainable AI

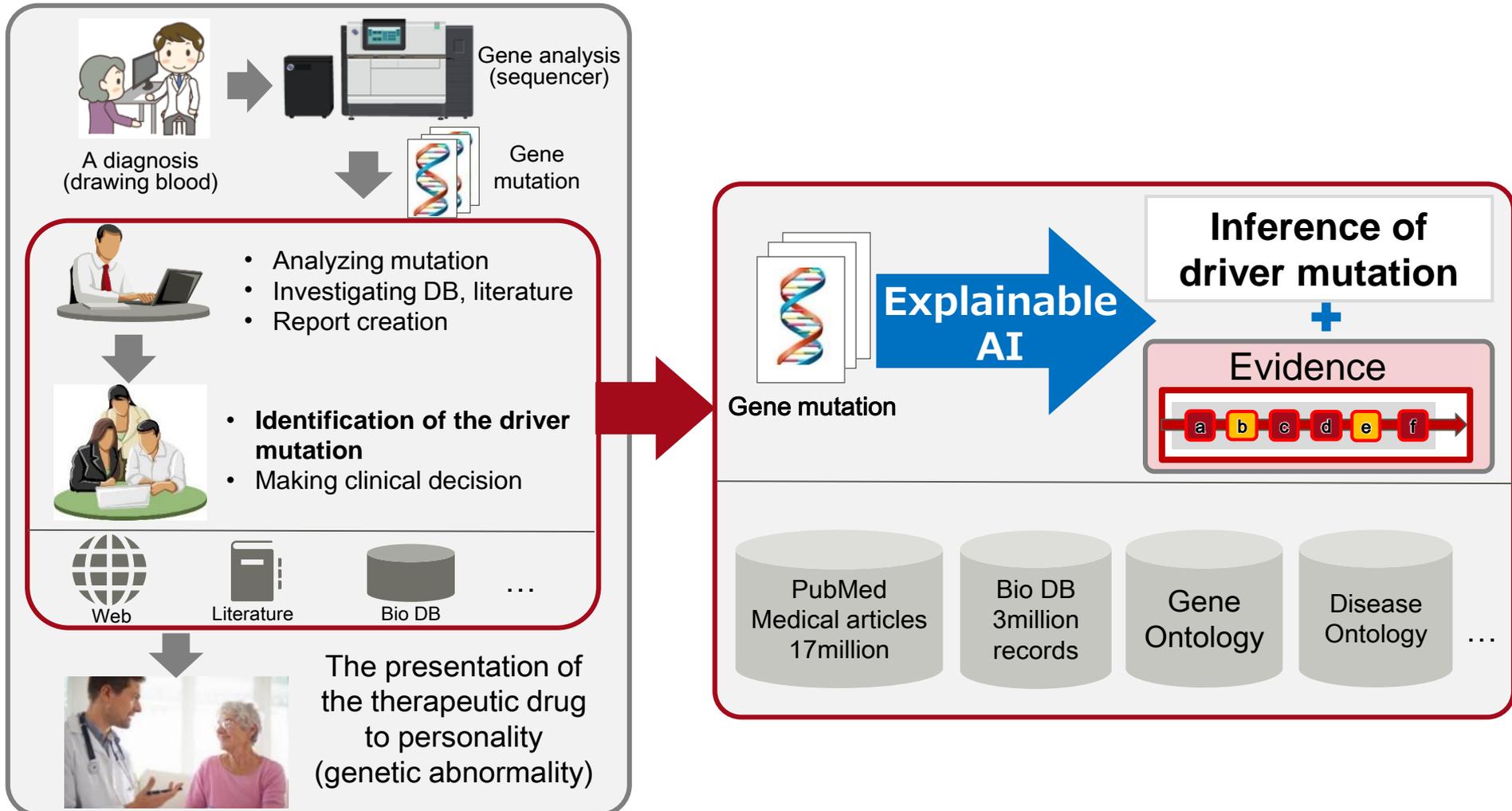- ■ If we applied **IN**-explainable AI to clinical decision.

She has cancer.
score = 0.99

She looks well.
**I can't believe AI.**
Let's see a state
for a while.

xx month ago

I found a small
cancer in her
stomach!

- ■ If we applied **explainable AI** to clinical decision.

**Explainable AI**

Output

Inference
result

Clinical
evidence

**Verify**

Cancer?
really?

**Cooperation** based on
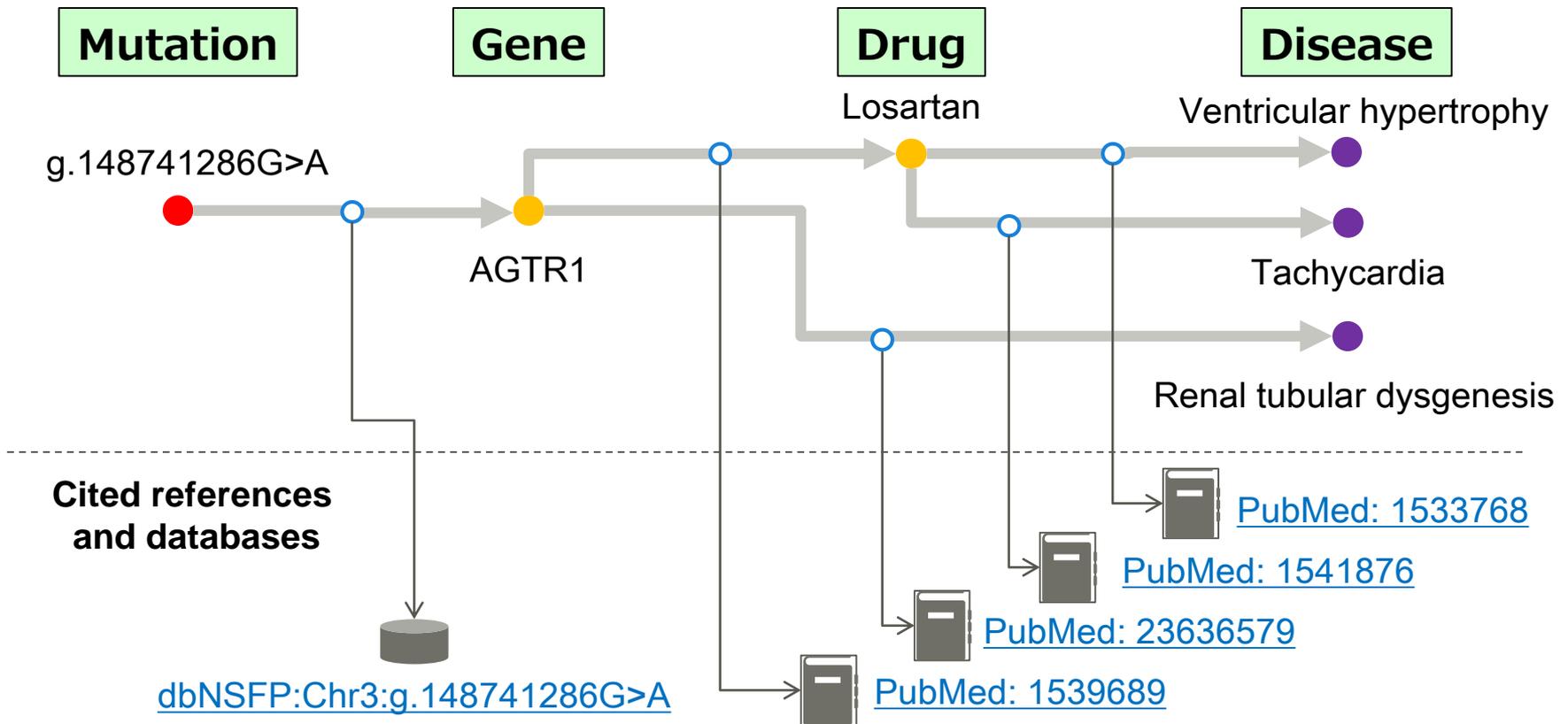the evidence report.

Let's carry out more
detailed examination

# Application to Genomic Medicine

■ We developed a prototype of AI that infers which gene mutation of the patient is associated with disease.



A diagnosis (drawing blood)

Gene analysis (sequencer)

Gene mutation

- Analyzing mutation
- Investigating DB, literature
- Report creation

- **Identification of the driver mutation**
- Making clinical decision

Web       Literature       Bio DB       …

The presentation of the therapeutic drug to personality (genetic abnormality)

Gene mutation

**Explainable AI**

**Inference of driver mutation**

+

Evidence

a b c d e f

| PubMed Medical articles 17million | Bio DB 3million records | Gene Ontology | Disease Ontology | … |

Ref. Hokkaido Univ. Hospital(http://www.huhp.hokudai.ac.jp/hotnews/detail/00001144.html)

# Formation of Evidence Path by Knowledge Graph

# Overview

## ■ Two key components

- **Deep Tensor** is a neural classifier of **graph data**

- **Knowledge graph** is an extremely large graph knowledge base

**(1) Explaining the important factor for the inference**

Outputting the factors which strongly influenced the inference result through Deep Tensor



**(2) Explaining the evidence for inference result**

Knowledge graph generates a evidence path from input to the inference result based on the inference factors.

# Knowledge Graph

# What is knowledge graph?

- ## Knowledge
  - Nodes and a edge that connects the nodes.

  subject          relation          object

  Francis Harry Compton Crick  $\xrightarrow{\text{found}}$  DNA double helix
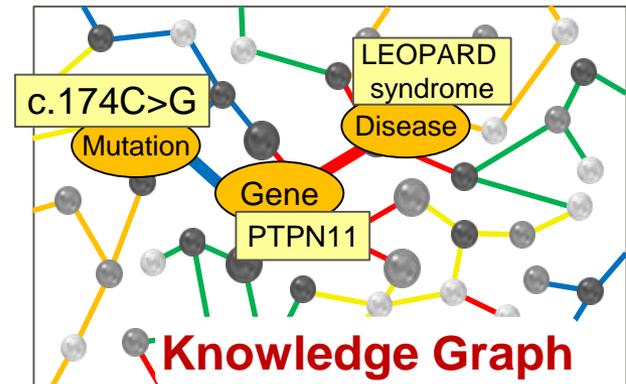
- ## Knowledge graph
  - Nodes
    - Gene, mutation, drug, disease, etc...
  - Edges
    - Drug A is responsive to disease B
    - Gene C has an function D
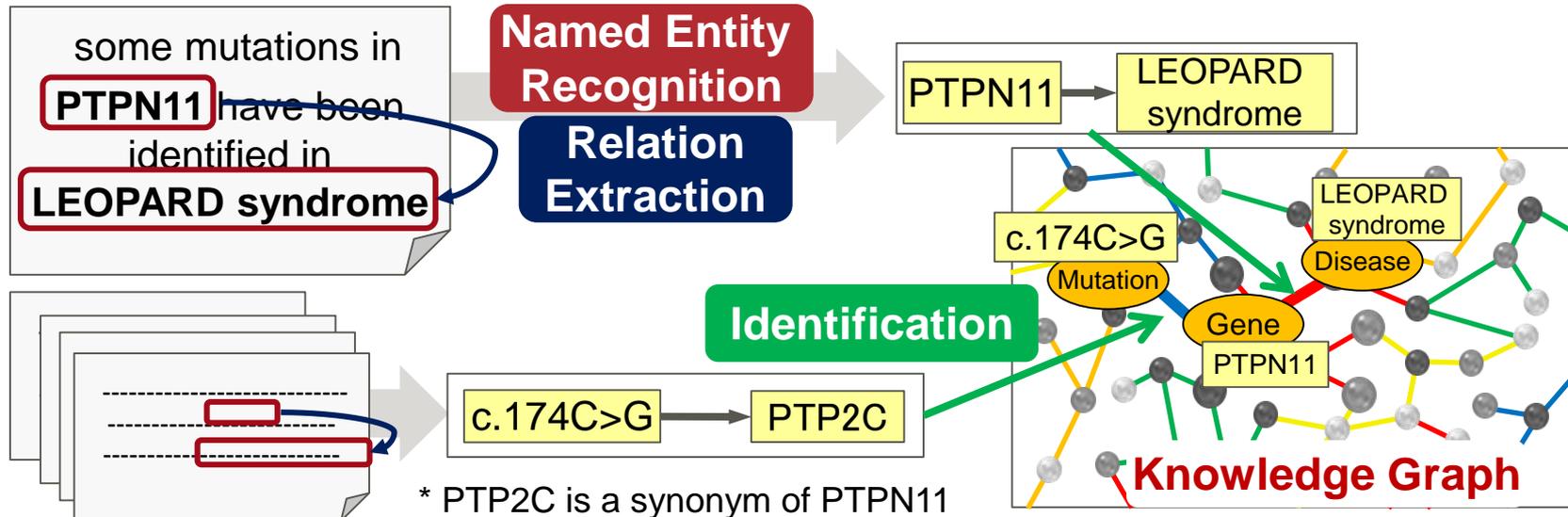    - Mutation E is located on Gene F
    - ...



c.174C>G
Mutation
LEOPARD syndrome
Disease
Gene
PTPN11

**Knowledge Graph**

# How to make knowledge graph

- **From existing database**
  - Extraction of relational data from public databases
    - Relations between gene and its attributes (name, ID, function, etc.)
    - Relations between protains

- **From literature**
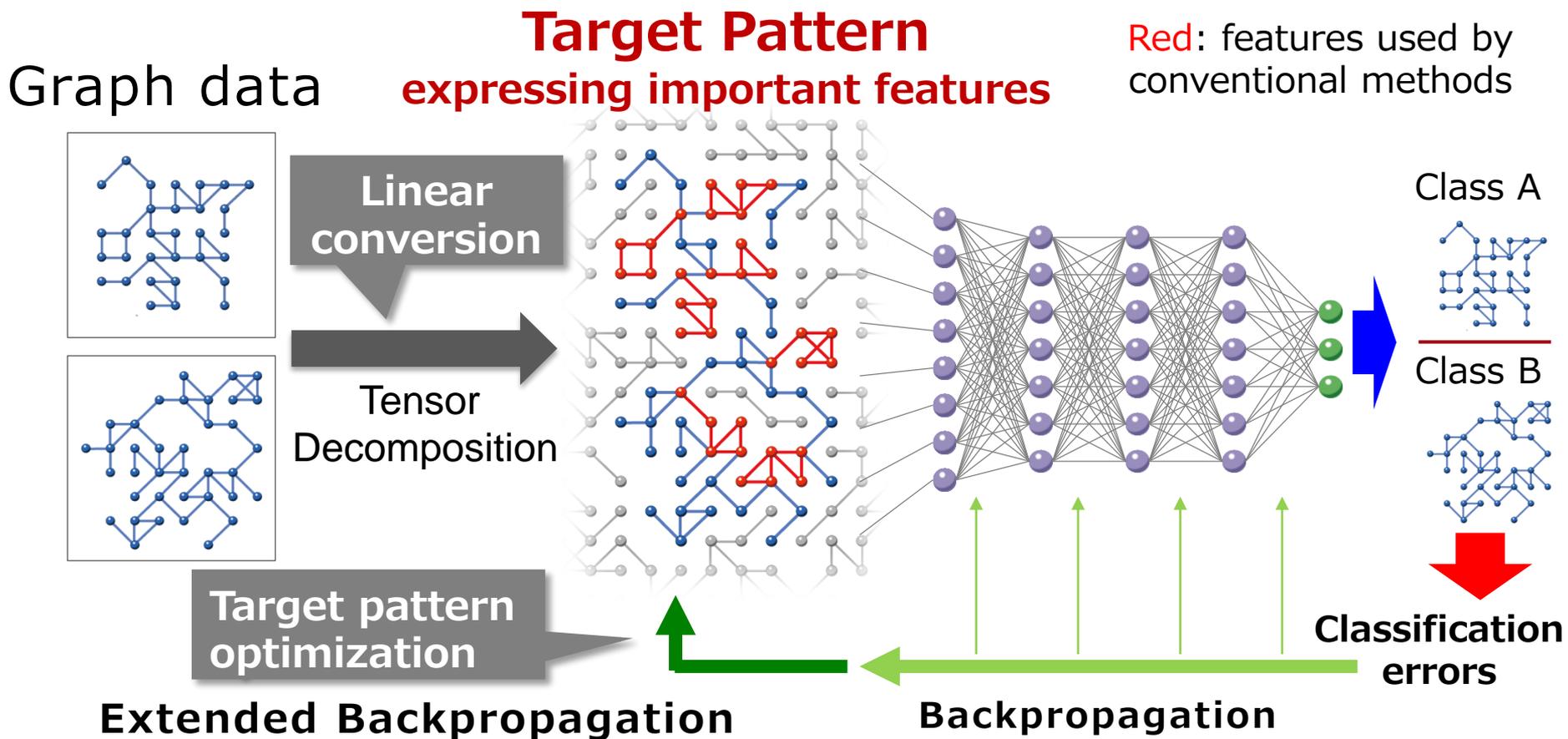  - Knowledge extraction using Natural Language Processing technologies



some mutations in **PTPN11** have been identified in **LEOPARD syndrome**

**Named Entity Recognition**

**Relation Extraction**

PTPN11 → LEOPARD syndrome

**Identification**

c.174C>G → PTP2C

\* PTP2C is a synonym of PTPN11

c.174C>G (Mutation)

LEOPARD syndrome (Disease)

Gene PTPN11

**Knowledge Graph**

# Deep Tensor

■ Koji Maruhashi, Masaru Todoriki, Takuya Ohwa, Keisuke Goto, Yu Hasegawa, Hiroya Inakoshi, Hirokazu Anai, Learning Multi-way Relations via Tensor Decomposition with Neural Networks, AAAI 2018 ([AAAI 2018](#)), February 2018.

# Naïve Idea

- **Linearly convert** graph data as close as a <u>target pattern</u>
- Optimize the target pattern along with Neural Networks

**Target Pattern**
**expressing important features**

Red: features used by conventional methods

Graph data



**Linear conversion**

Tensor Decomposition

**Target pattern optimization**

Class A

Class B

**Classification errors**

**Extended Backpropagation**

**Backpropagation**

## Classify graph data accurately with high interpretability

# Difficulties Using Graph for Deep Leaning

# Tensor Representation

■ Multi-way data can be represented as a tensor

**Multi-way data**

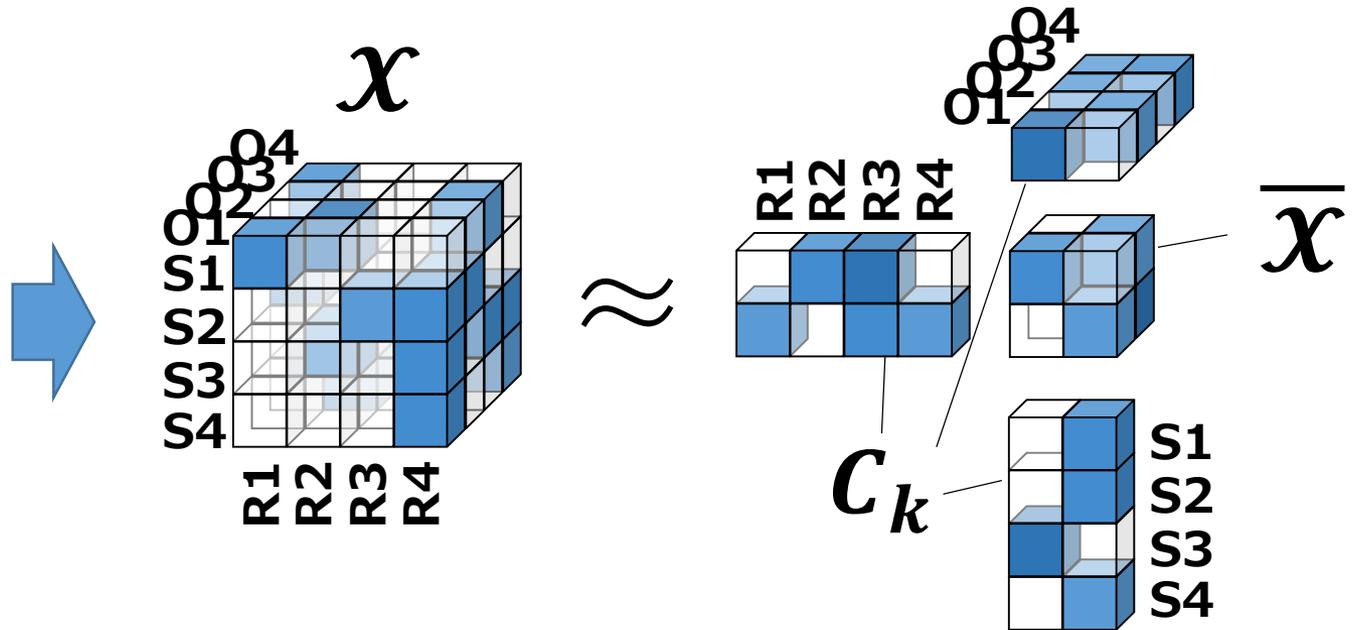| Sub | Relation | Obj |
|-----|----------|-----|
| S1 | R1 | O1 |
| S2 | R3 | O1 |
| S3 | R2 | O2 |
| ... | ... | ... |

**Tensor representation**

# Tensor Decomposition

- Approximate a tensor $\mathcal{X}$ by a core tensor $\overline{\mathcal{X}}$ multiplied by factor matrices $\{C_k\}$

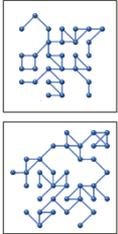- Results are <u>easy to interpret</u> in terms of nodes and edges.



$$(\overline{\mathcal{X}}, \{C_k\}) = \operatorname{argmin} \left\| \mathcal{X} - \overline{\mathcal{X}} \prod_k \times_k C_k^T \right\|_2^2$$

Can we leverage tensor decomposition?

# Leveraging Tensor Decomposition

■ Analyze graph data more efficiently

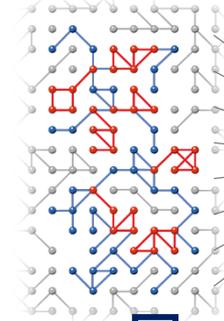Neural Network (NN)



Graph data

**Target Pattern**

Graph isomorphism determination

Too much cost!
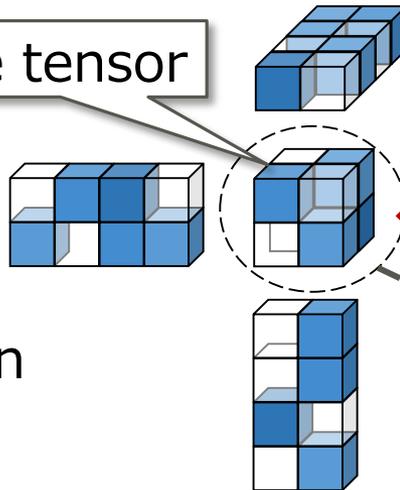
**Structure Restricted Tensor Decompsition**

**Tensor-based expression**

Core tensor

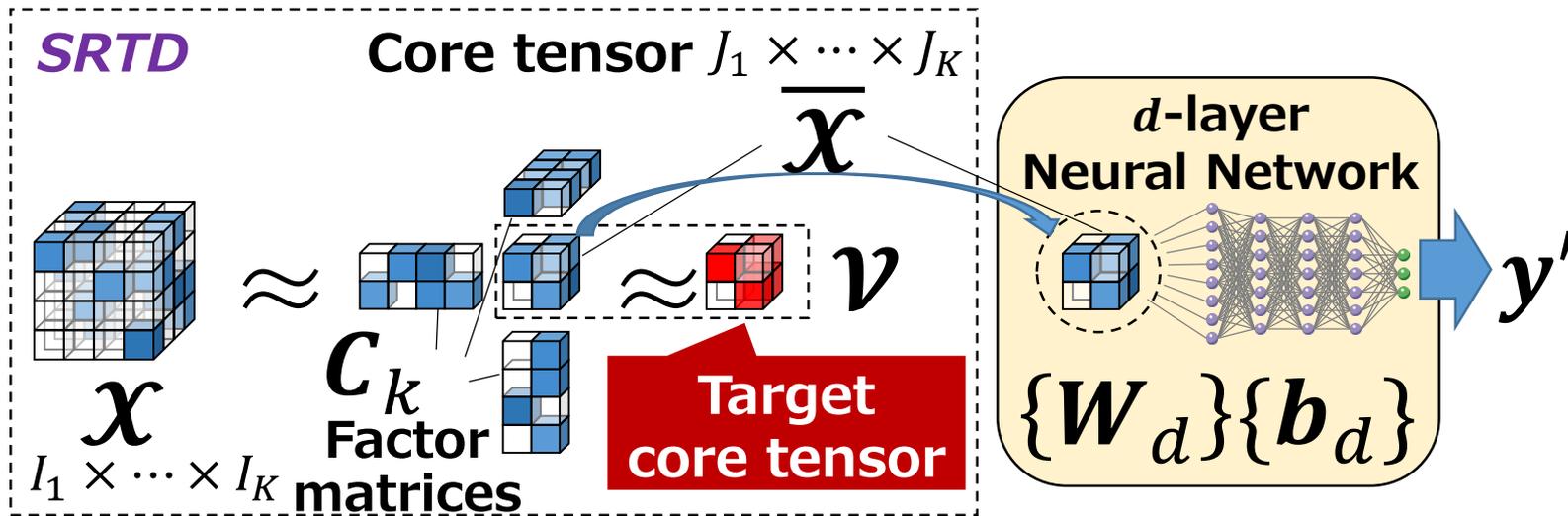**Target core tensor**

close

Optimize using Extended Backpropagation

≈

approximation

Input of NN

**A new type of tensor decomposition**

# Structure Restricted Tensor Decomposition (SRTD)

**Details**

- Calculate $\overline{\mathcal{X}}$ which minimize $\left\| \mathcal{X} \prod_{k|I_k < J_k} \times_k C_k - \overline{\mathcal{X}} \prod_{k|I_k \geq J_k} \times_k C_k^T \right\|_2^2$

- by using $C_k$ which minimize $\left\| \mathcal{X} \prod_{k|I_k < J_k} \times_k C_k - \mathcal{V} \prod_{k|I_k \geq J_k} \times_k C_k^T \right\|_2^2$

- subject to $C_k^T C_k = I (I_k \geq J_k), C_k C_k^T = I (I_k < J_k)$



*SRTD*   **Core tensor** $J_1 \times \cdots \times J_K$

$\overline{\mathcal{X}}$

$d$-**layer Neural Network**

$\mathcal{X} \approx$

$C_k$  **Factor matrices**

$\approx$  **Target core tensor**  $\mathcal{V}$

$I_1 \times \cdots \times I_K$

$\{W_d\}\{b_d\}$

$y'$

# How to interpret prediction results?

■ Solution: Learn _interpretable models_ that output similar results as Neural Networks



Tensor Decomposition is interpretable!

$\mathcal{X}$

$\overline{\mathcal{X}}$

Neural Network

**Black Box**

$\mathbf{c}_k$

$\mathcal{V}$

$y'$

**Similar results**

**Interpretable Model (Linear regression)**

$\overline{\mathcal{X}}$

**White Box**

$y''$

$$y'' = \langle \overline{\mathcal{X}}, \mathcal{W} \rangle + b$$

# Local Interpretable Models

- Use perturbed samples $\mathcal{X}^{(p)}$ based on LIME [Ribeiro+ KDD16]

  - Learn interpretable model $g_A\left(\overline{\mathcal{X}^{(p)}}\right) = \langle \overline{\mathcal{X}^{(p)}}, \mathbf{w}_A \rangle + \mathbf{b}_A$

  - which minimize $\sum_p \pi(p) \left\| y_A^{\prime(p)} - g_A\left(\overline{\mathcal{X}^{(p)}}\right) \right\|_2^2$

  - where $\pi(p) = \exp(- \left\| \overline{\mathcal{X}} - \overline{\mathcal{X}^{(p)}} \right\|_2^2 / \sigma^2)$



**Probability that $\mathcal{X}^{(p)}$ is classified into A by DeepTensor**

$$y_A^{\prime(1)} = 0.9 \quad y_A^{\prime(4)} = 0.1$$
$$y_A^{\prime(2)} = 1.0 \quad y_A^{\prime(5)} = 0.0$$
$$y_A^{\prime(3)} = 0.9 \quad y_A^{\prime(6)} = 0.1$$

- *Contribution score* is calculated by $\mathcal{X} * \mathbf{w}_A \prod_k \times_k \mathbf{c}_k^T$
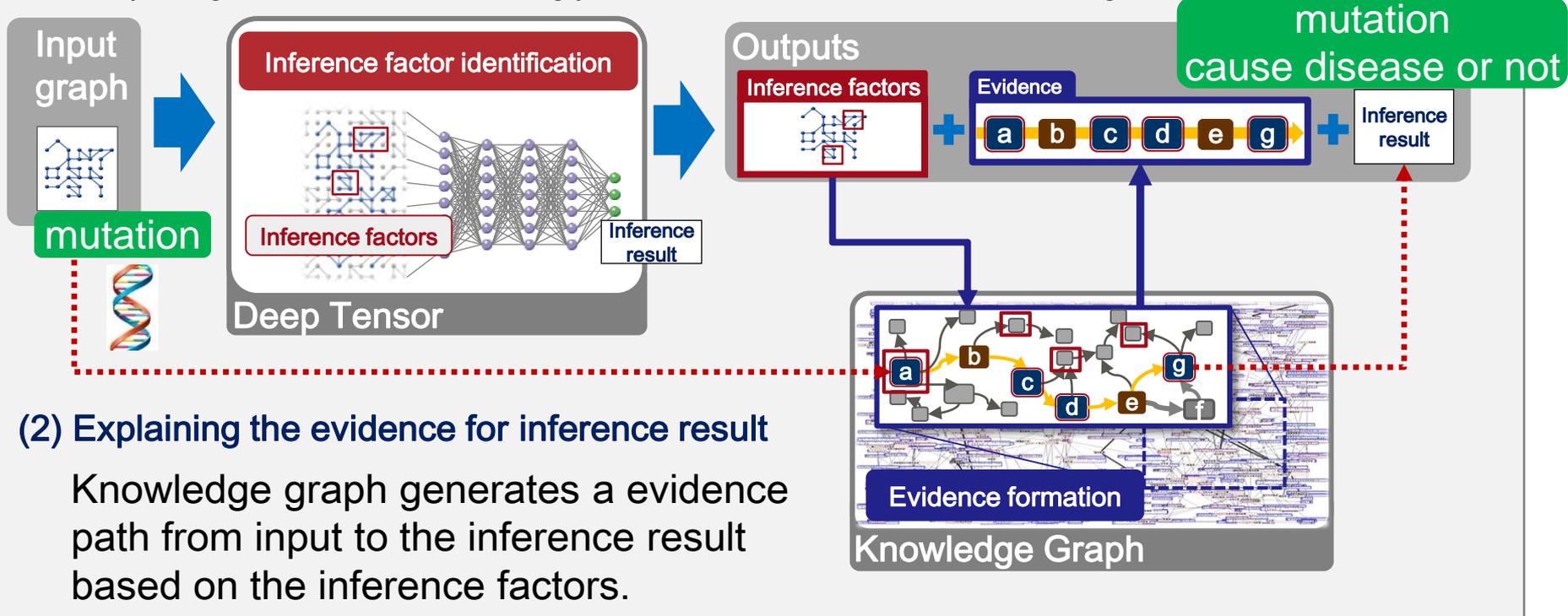
# Recap: Overview

**FUJITSU**

■ **Deep Tensor + knowledge graph**

■ Knowledge graph generates input graph (extracting subgraph representing about the mutation)

■ Deep Tensor infers which the mutation cause disease or not, and output inference factors.

■ Knowledge graph makes evidence graph based on the inference factors.

**(1) Explaining the important factor for the inference**
Outputting the factors which strongly influenced the inference result through Deep Tensor



**(2) Explaining the evidence for inference result**

Knowledge graph generates a evidence path from input to the inference result based on the inference factors.

# Conclusion

- Explainable AI is a key technology for cooperation of AI and humans

- We developed a prototype of explainable AI
  - The explainable AI explains important part of input data, and the evidence that explains the important part and inference result.

- We are now trying to proof of the concept of explainable AI, by cooperating with several medical groups.