

Causal Explanation

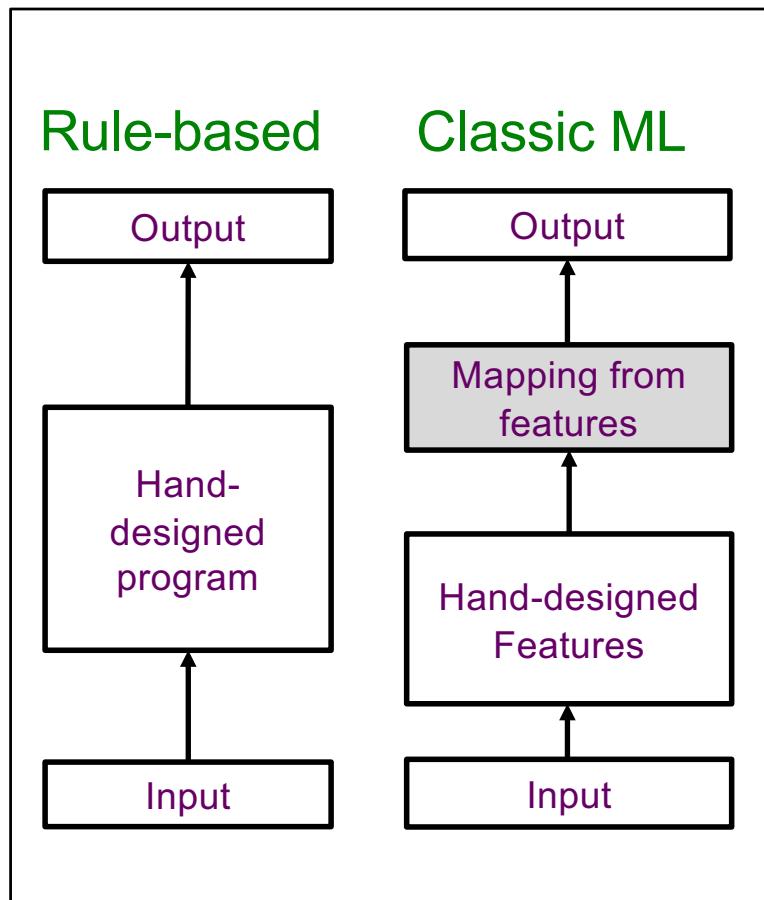
Sargur N. Srihari
srihari@buffalo.edu

Topics in Causal Explanation

- Three waves of AI
- Explainable AI
- Probabilistic Graphical Models
- Causal Modeling
- Most Probable Explanation
- Explaining Forensic Evidence

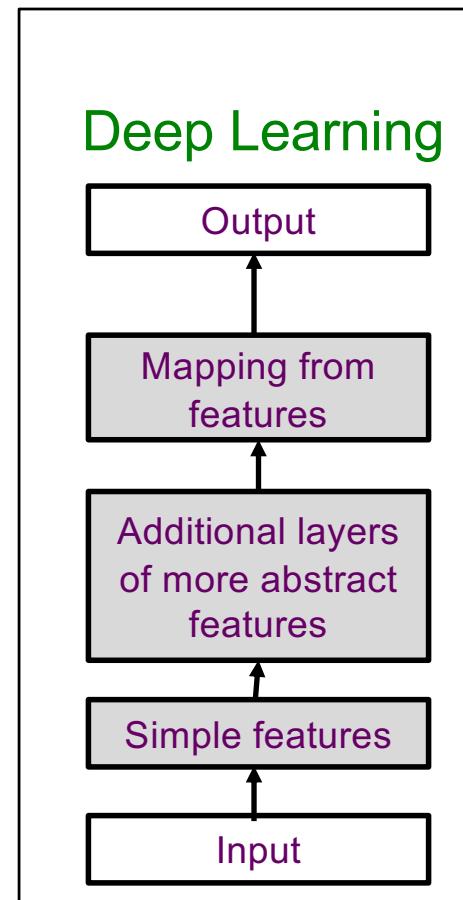
Role of probability in explanation

FIRST WAVE

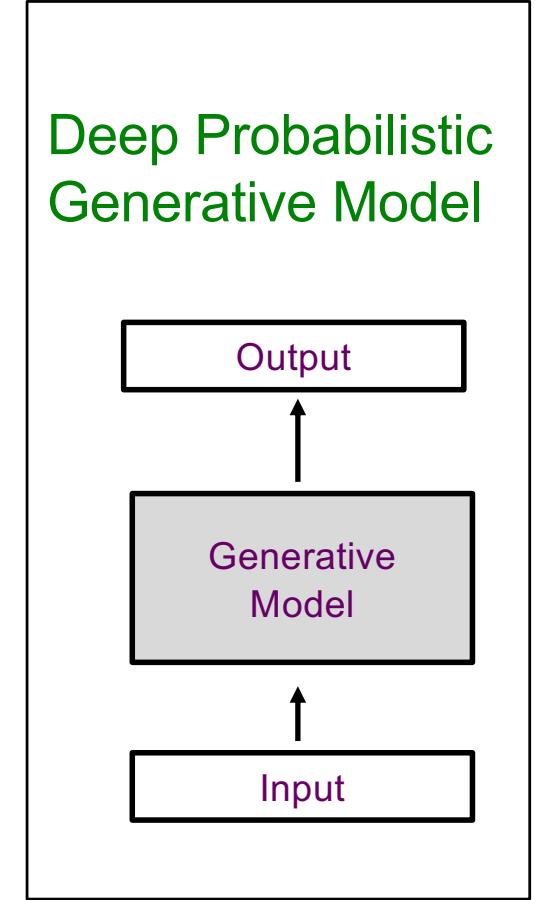


Naturally suited
for XAI

SECOND WAVE

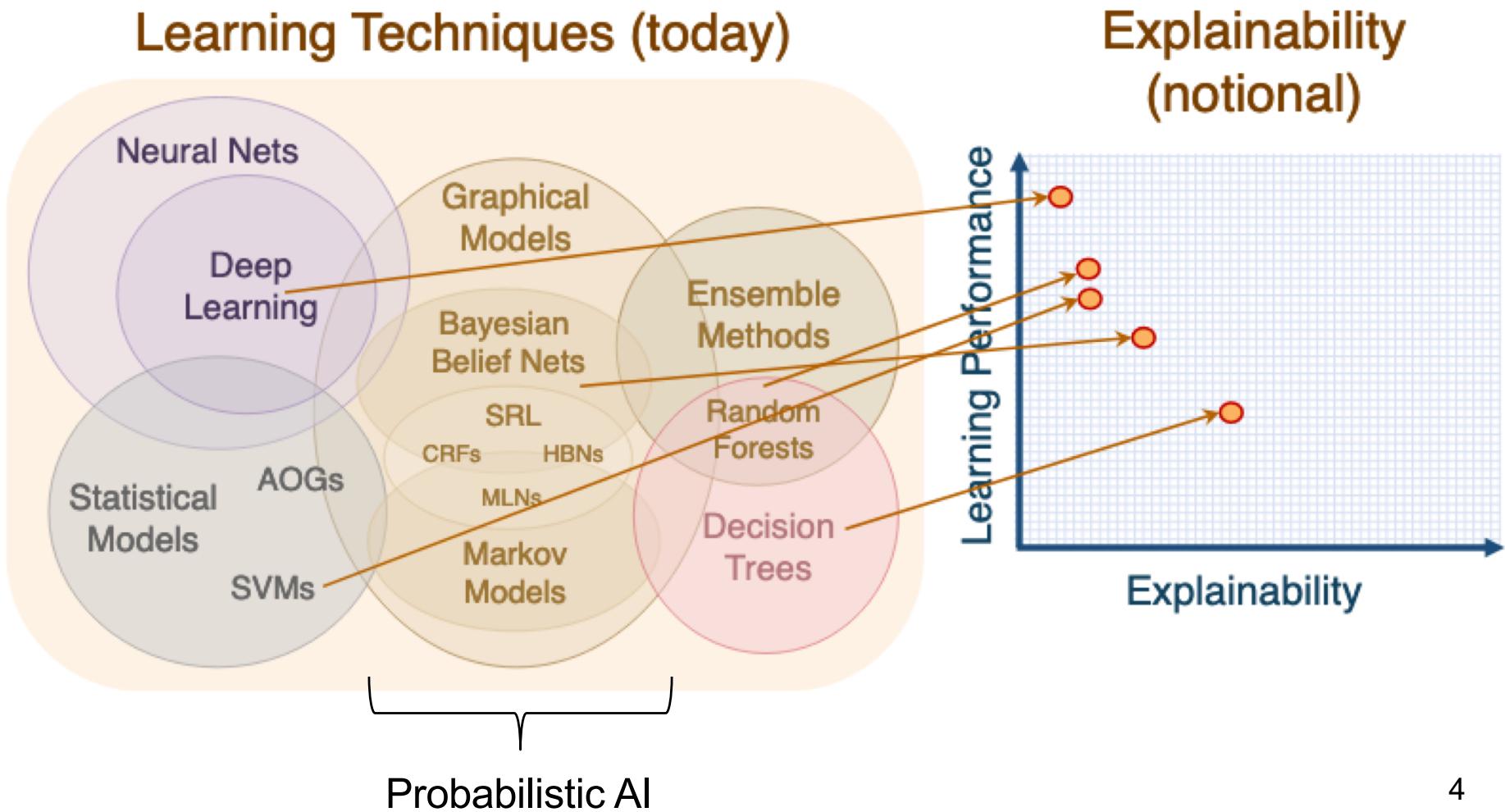


THIRD WAVE



Most Probable
Explanation

Explainability of AI Models

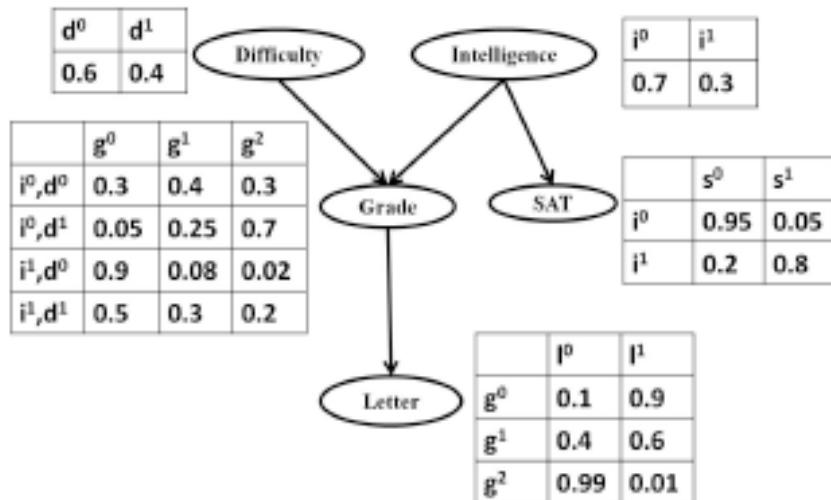


Probabilistic Models in Explanation

- Generative models are useful for
 - Querying
 - BNs, MRFs
 - Sampling
 - VAEs, GANs
- Inference methods useful as explanation
 - Most Probable Explanation (MPE)
 - Causal Inference
 - Explaining Away

Probabilistic Graphical Models (PGMs)

Bayesian Networks

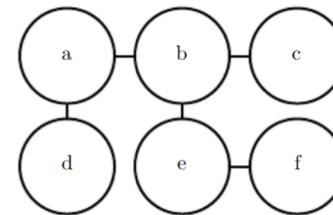


- **CPDs:** $p(x_i \mid pa(x_i))$
- **Joint Distribution** $\mathbf{x} = \{x_1, \dots x_n\}$

$$p(\mathbf{x}) = \prod_{i=1}^N p(x_i \mid pa(x_i))$$

$$\begin{aligned} P(D, I, G, S, L) &= \\ P(D)P(I)P(G \mid D, I)P(S \mid I)P(L \mid G) \end{aligned}$$

Markov Networks



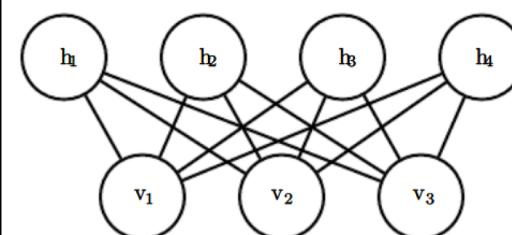
$$p(\mathbf{x}) = \frac{1}{Z} \hat{p}(\mathbf{x}) \quad \tilde{p}(\mathbf{x}) = \prod_{C \in G} \phi(C)$$

$$Z = \int \tilde{p}(\mathbf{x}) d\mathbf{x}$$

$$p(a, b, c, d, e, f) = \frac{1}{Z} \phi_{a,b}(a, b) \phi_{b,c}(b, c) \phi_{a,d}(a, d) \phi_{b,e}(b, e) \phi_{e,f}(e, f)$$

Energy model $E(a, b, c, d, e, f) =$
 $E_{a,b}(a, b) + E_{b,c}(b, c) + E_{a,d}(a, d) + E_{b,e}(b, e) + E_{e,f}(e, f)$
 $\phi_{a,b}(a, b) = \exp(-E(a, b))$

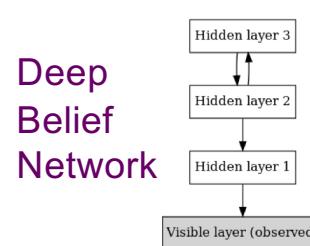
Restricted Boltzmann machine



$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^T \mathbf{v} - \mathbf{c}^T \mathbf{h} - \mathbf{v}^T \mathbf{W} \mathbf{h}$$

$$p(\mathbf{h}|\mathbf{v}) = \prod_i p(h_i|v)$$
 and

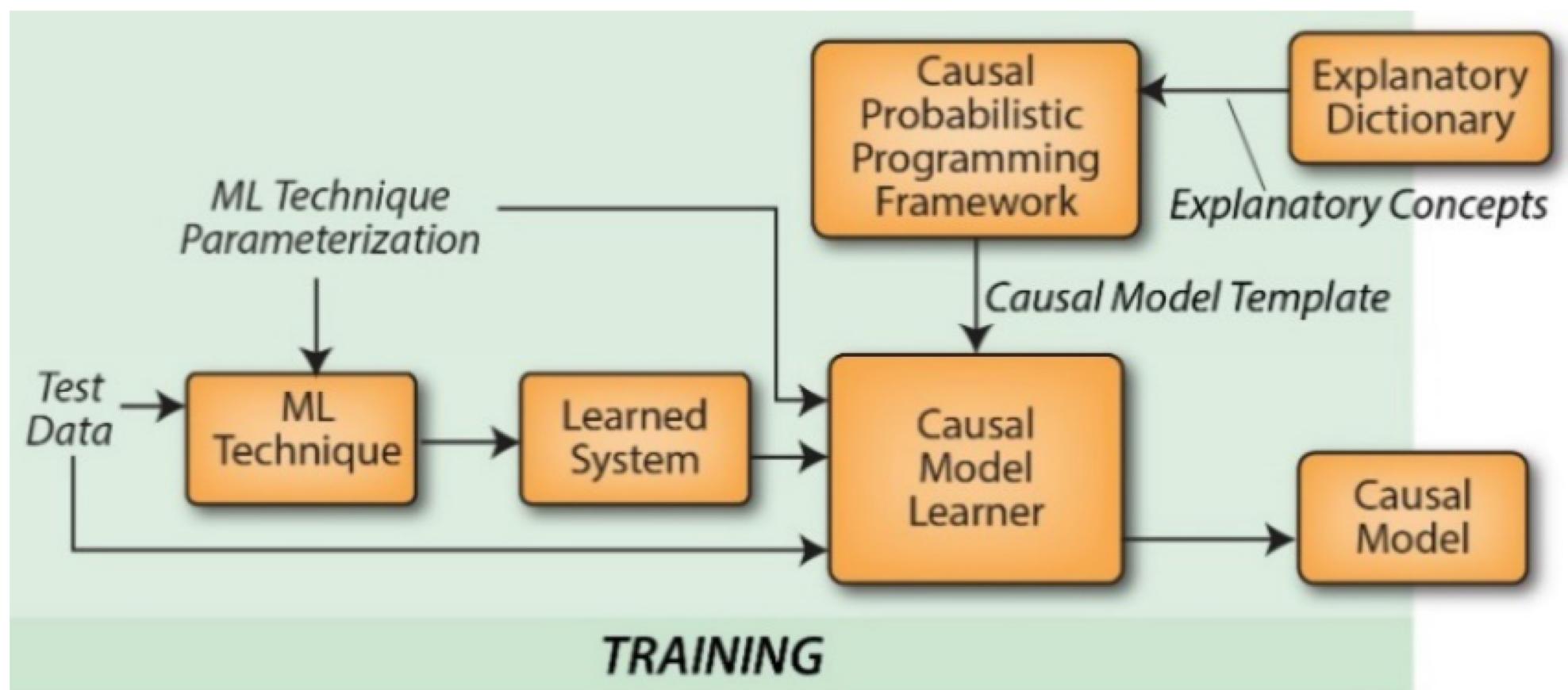
$$p(\mathbf{v}|\mathbf{h}) = \prod_i p(v_i|h)$$



Deep
Belief
Network

XAI from Causal Model Induction

- Experiment with the learned model to learn an explainable, causal, probabilistic programming model



Causal Modeling

- For a G defining a causal directed PGM on variables $X=x_1, \dots, x_n$, the joint distribution is

$$P(x_1, \dots, x_n) = \prod_i P(x_i | pa_i)$$

- where pa_i are the parents of x_i
- An intervention on G by setting $x_i = x'_i$, denoted as $do(x'_i)$, induces a modified graph G' where the edges from pa_i to x_i are removed, resulting in post intervention distribution

$$P(x_1, \dots, x_n | do(x'_i)) = \begin{cases} \prod_{j \neq i} P(x_j | pa_j) & \text{if } x_i = x'_i \\ 0 & \text{if } x_i \neq x'_i \end{cases}$$

- Semantics of causality are important for explaining DNNs, because explanations must be causal models

Explanation of DNN decision

- When one seeks an explanation of a network's decisions, one is equivalently asking "What changes can be made to the input for the output to change or stay the same?"
- Consider a causal model that defines the joint distribution $P(\emptyset, \mathbb{P}, \mathbb{X})$ over a set of DNN outputs \emptyset , inputs \mathbb{P} , & intermediate variables \mathbb{X}
- Explanation of an observed DNN output is formulated as intervention

Causal Representation in DNNs

- A user could ask counterfactual questions about the network, i.e. $P(\emptyset, \mathbb{P}, \mathbb{X} \mid \text{do}(x'_i))$ for any input, output, or internal neuron in the network
- But it serves poorly as a model of explanation
 - due to the lack of human-level concepts that underlie any arbitrary neuron in the network:
 - saying neuron i caused the network to detect a pedestrian may be correct but does not satisfy human needing explanations
- Thus, a DNN causal model must be at a granularity meaningful to humans.

Mapping Neurons to Concepts

- A causal model for DNNs should be represented by joint distribution over \mathbb{O}, \mathbb{P} , and a set of concepts \mathbb{C}
- The process for deriving \mathbb{C} is described by a function $f_{\mathbb{R}}: \mathbb{R} \rightarrow \mathbb{C}$ over a specific DNN that transforms the representation of neurons and their activations into a set of concept variables
- Ideally, $f_{\mathbb{R}}$ would have the following properties:

$$\int_R P(\mathbb{O}, \mathbb{P}, \mathbb{R}) = \int_C P(\mathbb{O}, \mathbb{P}, \mathbb{C})$$

$$P(\mathbb{O}, \mathbb{P}|R, do(p'_i)) = P(\mathbb{O}, \mathbb{P}|C, do(p'_i))$$

Concept Extraction

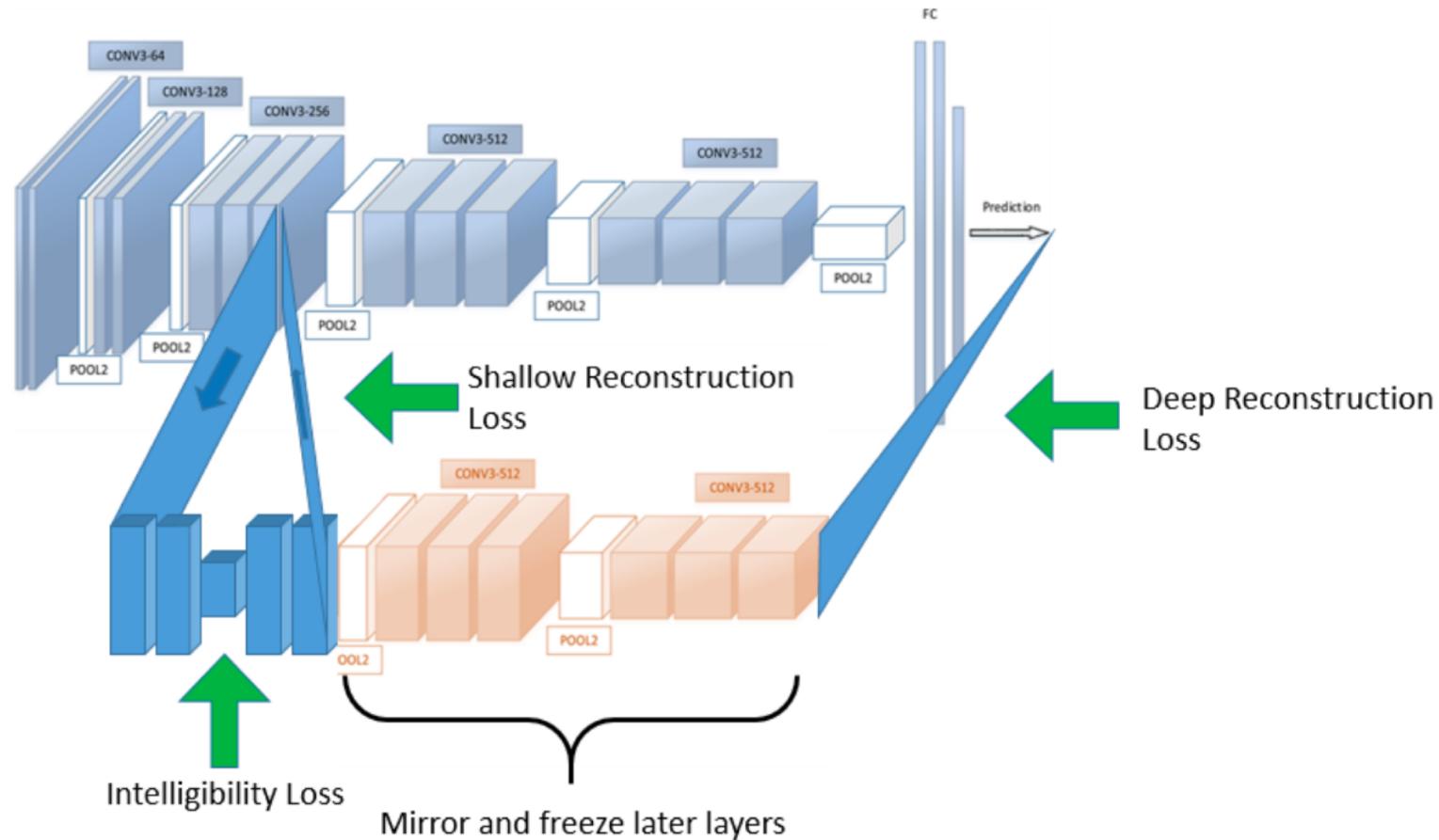
- One way to construct concepts that satisfy the semantics is to consider network activations
- But concepts needn't be restricted to those
 1. Concepts should be few to minimize amount of investigation a human would need to employ
 2. Concepts should be interpretable
 - in the case of images, we would like activations to be restricted to contiguous areas containing consistent, interpretable visual features
 3. Concepts should contain all relevant information needed for achieving network task

Auxiliary neural network

- Network constructs concept representations satisfying above properties
- Using an autoencoder with two loss functions
 - Shallow reconstruction loss: L_1 norm between input and output activations
$$L_{shallow}(\theta; a_i) = |d_\theta(c_\theta(a_i)) - a_i|_1$$
$$L_{deep}(\theta; a_i) = KL(r(a_i) || r(d_\theta(c_\theta(a_i))))$$
 - Deep reconstruction loss: reconstructed activations result in same classification output after being passed through the rest of network
 - Total autoencoder loss

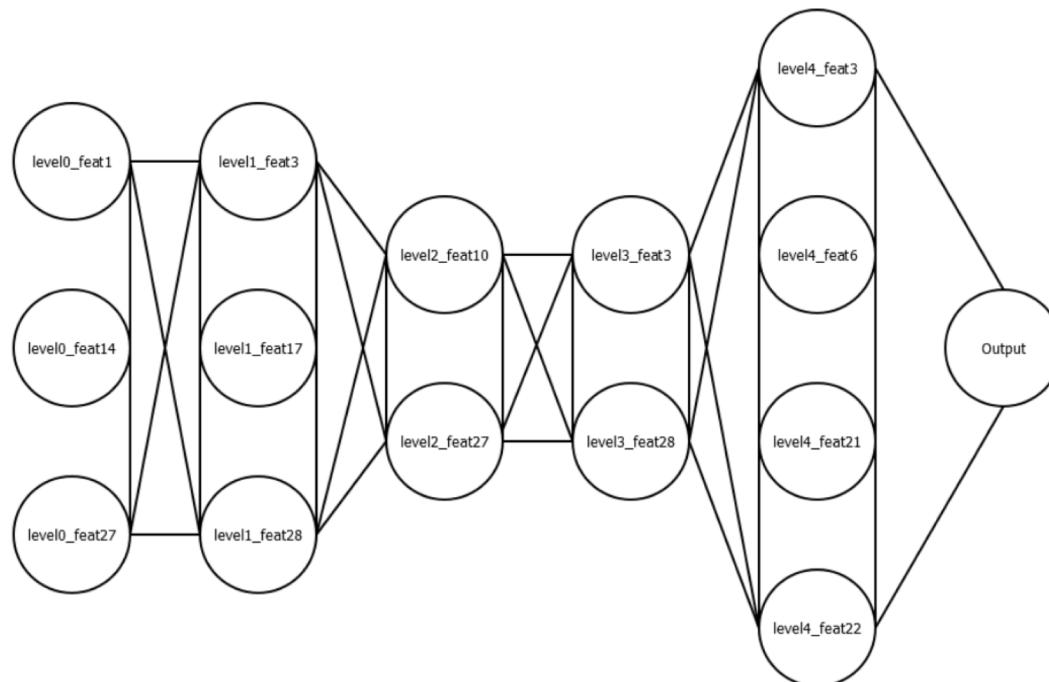
$$L(\theta; x_i) = \lambda_{shallow} L_{shallow}(\theta; x_i) +$$
$$\lambda_{deep} L_{deep}(\theta; x_i) +$$
$$\lambda_{interpretability} L_{interpretability}(\theta; x_i)$$

Autoencoder for a single layer



Learned causal Bayes Net

- For Inria pedestrian data set



Crossed boxes indicate that all nodes of a given level have edges incident on each of the nodes of the subsequent level.

Summary of Causal Learning XAI

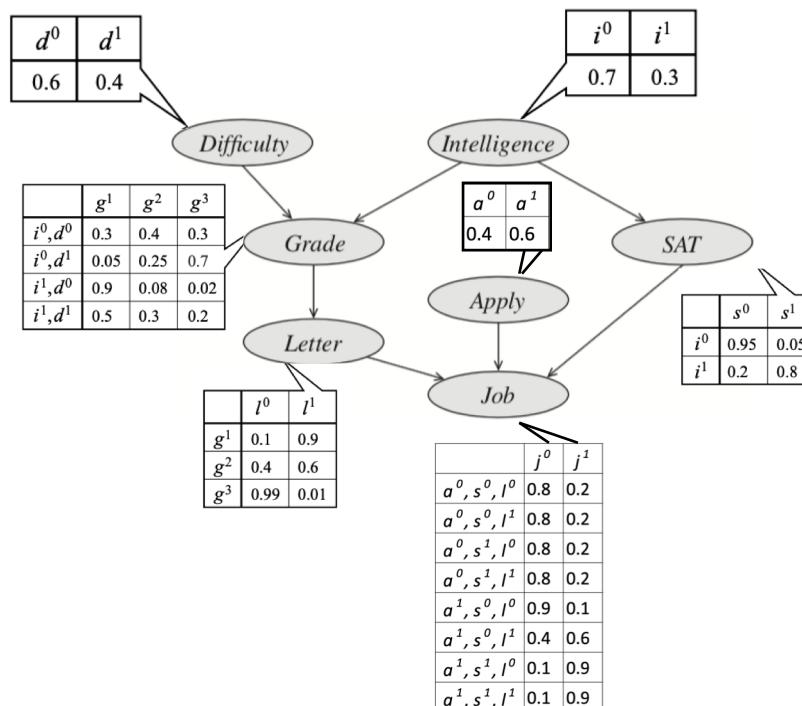
- Approach to explain predictions to relate output prediction to concepts represented within
 - Series of autoencoders encouraging interpretable properties to construct concepts in activations throughout network
 - Pool features and intervene on autoencoded network to construct variables used to build a causal BN
 - Use network to identify features of causal relevance to individual classifications which are then visualized

Most Probable Explanation (MPE)

Observe evidence: student has received a job offer $J=j^1$.

We want an explanation for each of the following two queries

1. Most probable (MAP) setting for all variables I, D, G, S, L, A
2. Most probable setting for a particular variable I



$$P(I, D, G, S, L, A, J) = P(I)P(D)P(G|I, D)P(S|I)P(L|G)P(A)P(J|A, S, L)$$

Query 1: Inferring all variables

Conditioning with Job with $J=j^1$

Conditional Probability:

$$P(I, D, G, S, L, A | J = j^1) = \frac{P(I, D, G, S, L, A, J = j^1)}{P(J = j^1)}$$

Numerator:

$$P(I, D, G, S, L, A, J = j^1) = P(I)P(D)P(G|D, I)P(L|G)P(S|I)P(A)P(J = j^1|A, S, L)$$

Denominator requires Marginal Inference

$$P(J = j^1) = \sum_I \sum_D \sum_G \sum_S \sum_L \sum_A P(I, D, G, S, L, A, J = j^1) = 0.3679$$

MPE for all variables

$$P(I, D, G, S, L, A | J = j^1) = \frac{P(I, D, G, S, L, A, J=j^1)}{P(J=j^1)}$$

Exhaustive listing of 96 Combinations of numerator with probabilities

Sr. No.	Difficult	Intelligenc	Grade	Letter	SAT	Apply	Job	Joint Probabil
3	0	1	0	1	1	1	1	0.0629856
5	0	0	0	1	0	1	1	0.0387828
10	0	0	1	1	0	1	1	0.0344736
13	1	1	0	1	1	1	1	0.023328
16	1	0	2	0	0	0	1	0.01474704
18	1	0	1	1	0	1	1	0.014364
21	1	0	2	0	0	1	1	0.01106028
22	0	1	0	1	0	1	1	0.0104976
23	1	1	2	0	1	1	1	0.01026432
25	0	0	2	0	0	0	1	0.00948024
27	0	1	0	1	1	0	1	0.0093312
28	1	1	1	1	1	1	1	0.0093312
29	0	0	0	1	0	0	1	0.0086184
31	0	0	1	1	0	0	1	0.0076608
32	0	0	2	0	0	1	1	0.00711018
33	0	1	0	0	1	1	1	0.0069984
37	1	1	1	0	1	1	1	0.0062208

Max of numerator (a joint probability):

$$P(I, D, G, S, L, A, J=j^1) = 0.9 * 0.6 * 0.8 * 0.9 * 0.9 * 0.3 * 0.6 = 0.0629856$$

Max Conditional Probability (MAP):

$$P(I, D, G, S, L, A | J = 1) = 0.0629856 / 0.3679 = 0.1712$$

MPE is:

- $I = i^1$: Is Intelligent
- $D = d^0$: course is Difficult
- $G = g^0$: Grade is A
- $S = s^1$: high SAT
- $L = l^1$: strong Letter
- $A = a^1$: Apply for Job

Query 2: Marginal for I given $J = j^1$

The most probable explanation of *Intelligence* given $Job : J=j^1$

Answer: $I = i^1$: Is Intelligent

Formula:

$$P(I|J=1) = P(I, J)/P(J=1)$$

$$P(I, J) = P(I)P(J|A, S, L)P(A)P(S|I)P(L|G)P(G|D, I)P(D)$$

$$P(I=0, J=1) = 0.9 * 0.6 * 0.8 * 0.9 * 0.9 * 0.3 * 0.6 = 0.0629$$

$$P(I=1, J=1) = 0.6 * 0.7 * 0.3 * 0.9 * 0.95 * 0.6 * 0.6 = 0.0387$$

$$P(I=0|J=1) = 0.0387/(0.0629 + 0.0387) = 0.381$$

$$P(I=1|J=1) = 0.0629/(0.0629 + 0.0387) = 0.619$$

MAP of Student Intelligence

intel_0	0.381
intel_1	0.619

Compare to Prior of Student Intelligence



i^0	i^1
0.7	0.3

MPE is NP-Hard

- BN is an MRF whose factors are conditional probabilities and $Z=1$
- MAP inference is simpler than Marginal Inference since we can ignore Z

Marginal

$$p(y = 1) = \sum_{x_1} \sum_{x_2} \cdots \sum_{x_n} p(y = 1, x_1, x_2, \dots, x_n)$$

← Needs
Partition function

MAP

$$\max_{x_1, \dots, x_n} p(y = 1, x_1, \dots, x_n)$$

← Can ignore
Partition function

- However MAP inference is still not easy
 - Both Marginal and MAP inference are NP-hard

Sampling from Student BN

Algorithm 12.1 Forward Sampling in a Bayesian network

```
Procedure Forward-Sample (
     $\mathcal{B}$  // Bayesian network over  $\mathcal{X}$ 
)
1   Let  $X_1, \dots, X_n$  be a topological ordering of  $\mathcal{X}$ 
2   for  $i = 1, \dots, n$ 
3        $\mathbf{u}_i \leftarrow \mathbf{x} \langle \text{Pa}_{X_i} \rangle$  // Assignment to  $\text{Pa}_{X_i}$  in  $x_1, \dots, x_{i-1}$ 
4       Sample  $x_i$  from  $P(X_i | \mathbf{u}_i)$ 
5   return  $(x_1, \dots, x_n)$ 
```

PGMPY:

```
from pgmpy.sampling import BayesianModelSampling
forward_sampler = BayesianModelSampling(student)
samples = forward_sampler.forward_sample(size=1000)
```

	apply	intel	sat	diff	grade	letter	job
0	1	1	1	1	0	1	1
1	0	0	0	0	0	1	0
2	1	1	1	1	1	1	1
3	1	0	0	1	2	0	0
4	1	1	1	0	2	0	1
5	0	0	0	1	2	0	1
6	1	0	0	0	0	1	0
7	1	0	1	0	1	1	1
8	1	0	0	0	1	1	1
9	1	0	0	1	1	0	0
...
993	0	0	1	1	1	1	0
994	1	1	1	0	0	1	1
995	1	0	0	0	2	0	0
996	0	0	0	0	0	1	0
997	1	0	0	1	2	0	0
998	0	0	0	0	0	0	1
999	1	0	0	0	2	0	0

1000 rows × 7 columns

PGM classification over 1000 samples

```
from pgmpy.models import BayesianModel
from pgmpy.factors.discrete import TabularCPD
from pgmpy.inference import VariableElimination

student = BayesianModel()
student.add_nodes_from(['diff', 'intel', 'grade', 'sat', 'letter', 'apply', 'job'])
student.add_edges_from([('diff', 'grade'),
                      ('intel', 'grade'),
                      ('intel', 'sat'),
                      ('grade', 'letter'),
                      ('apply', 'job'),
                      ('letter', 'job'),
                      ('sat', 'job')])
student.fit(train)

inference = VariableElimination(student)
counter = 0
for index, row in test.iterrows():
    mle1 = inference.map_query(variables=['sat', 'letter'],
                                evidence={'job': row['job']})
    mle2 = inference.map_query(variables=['intel'],
                                evidence={'sat': mle1['sat'],
                                          'letter': mle1['letter']})
    if(mle2['intel'] == row['intel']):
        counter += 1

print('Accuracy is: ' + str(counter/test.shape[0]*100))

Accuracy is: 88.125
```

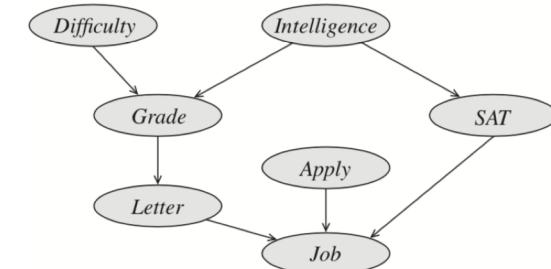


Figure 5.3 The Student example augmented with a *Job* variable

Results:

- Training Accuracy: 92%
- **Testing Accuracy: 88%**

Logistic Regression for Intelligence

Training Hyperparameters:

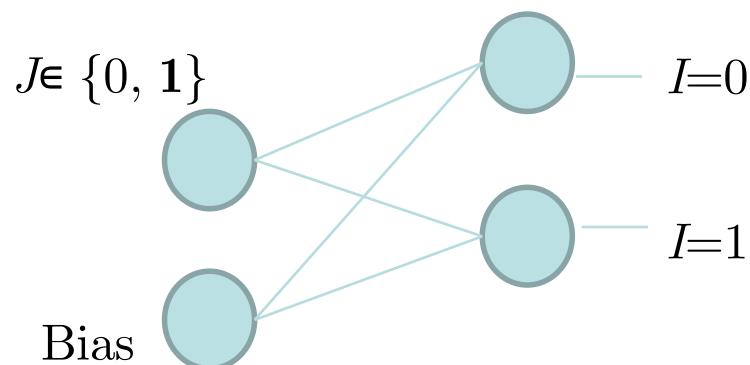
- Number of Epochs: 100
- Batch Size: 100
- Number of samples: 1000
- Number of Training samples: 800
- Number of Validation samples: 200

Results:

- Training Accuracy: 94.62%
- Training Loss: 0.1580
- **Testing Accuracy: 96%**
- **Testing Loss: 0.1534**

Logistic Regression output:

	Intelligence = 0	Intelligence = 1
Job = 1	0.038	0.962



Explanation:

MPE

evidence: student **gets job**

explanation: 96% probability that student is **intelligent**;

Formula:

$$P(I = 0|J = 1) = 0.498$$

$$P(I = 1|J = 1) = 0.502$$

Explainability vs Performance

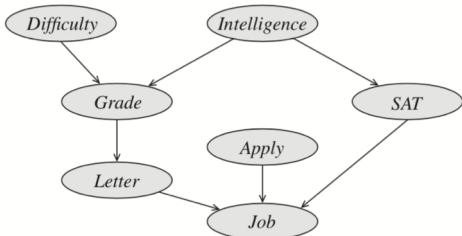
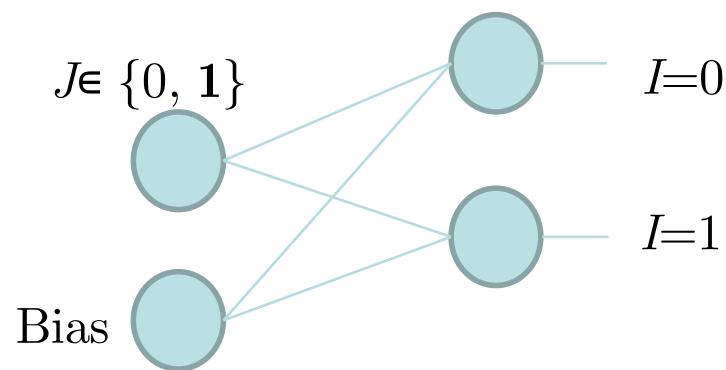


Figure 5.3 The Student example augmented with a *Job* variable

```
inference = VariableElimination(student)
counter = 0
for index, row in test.iterrows():
    mle1 = inference.map_query(variables=['sat', 'letter'],
                                evidence={'job': row['job']})
    mle2 = inference.map_query(variables=['intel'],
                                evidence={'sat': mle1['sat'],
                                          'letter': mle1['letter']})
    if(mle2['intel'] == row['intel']):
        counter += 1
print('Accuracy is: ' + str(counter/test.shape[0]*100))
Accuracy is: 88.125
```

Explainable System Results:

- Training Accuracy: 92%
- **Testing Accuracy: 88%**

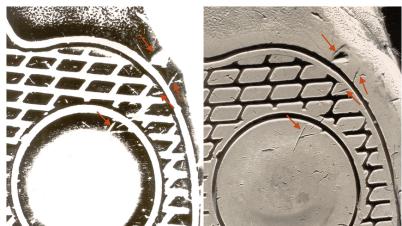


High Performance System Results:

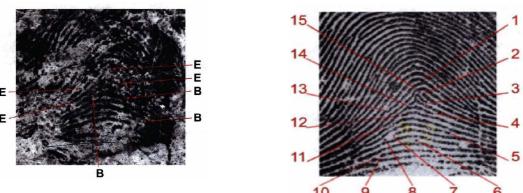
- Training Accuracy: 95%
- **Testing Accuracy: 96%**

Forensic Comparison

- Impression Evidence:
materials with characteristics of impressed objects
 - Footwear impressions



- latent fingerprints



- Handwriting samples



Expert Features

Feature Vector

imagename	0968c_num1.png
pen_pressure	2
letter_spacing	2
size	2
dimension	1
is_lowercase	2
is_continuous	2
slantness	3
tilt	2
entry_stroke_a	1
staff_of_a	2
formation_n	2
staff_of_d	3
exit_stroke_d	2
word_formation	2
constancy	1

Number of Features: 15

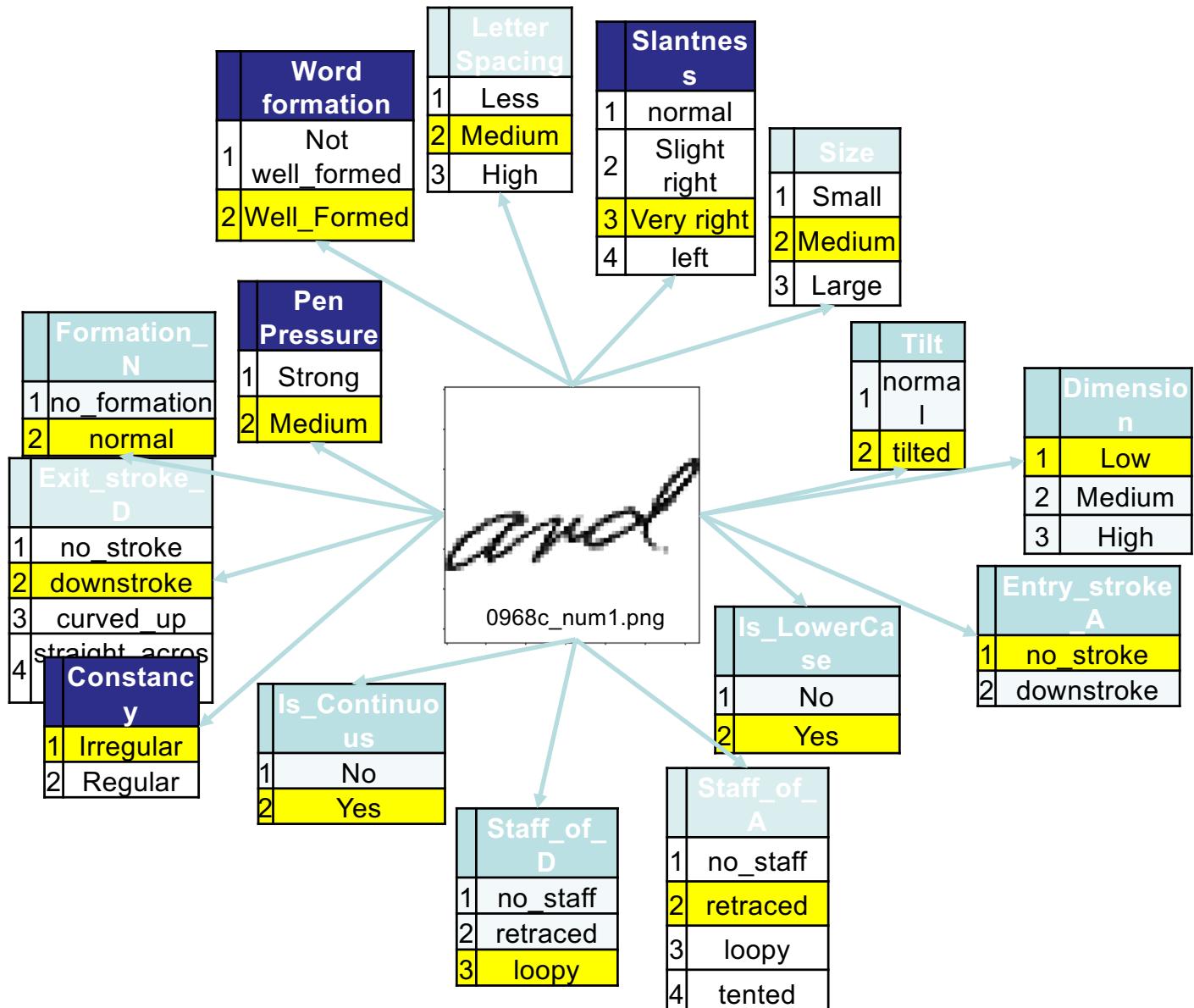
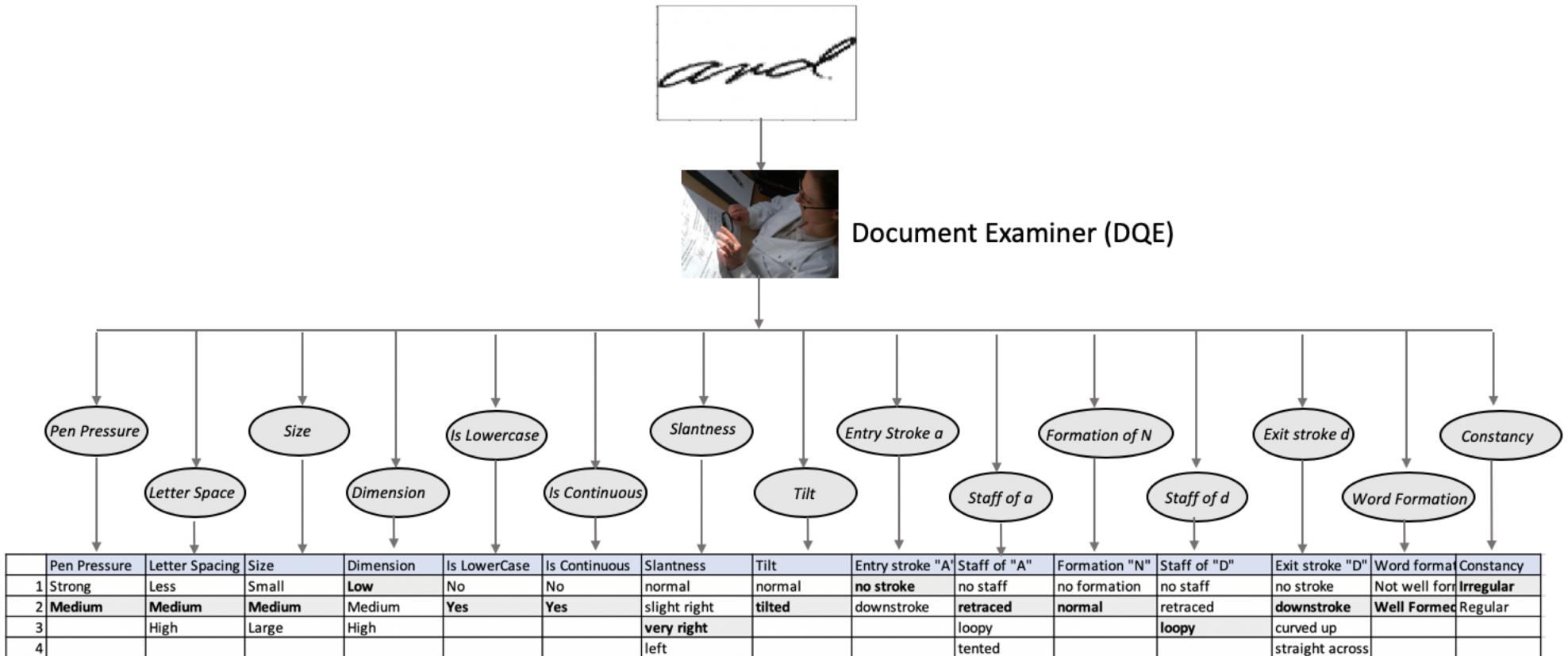
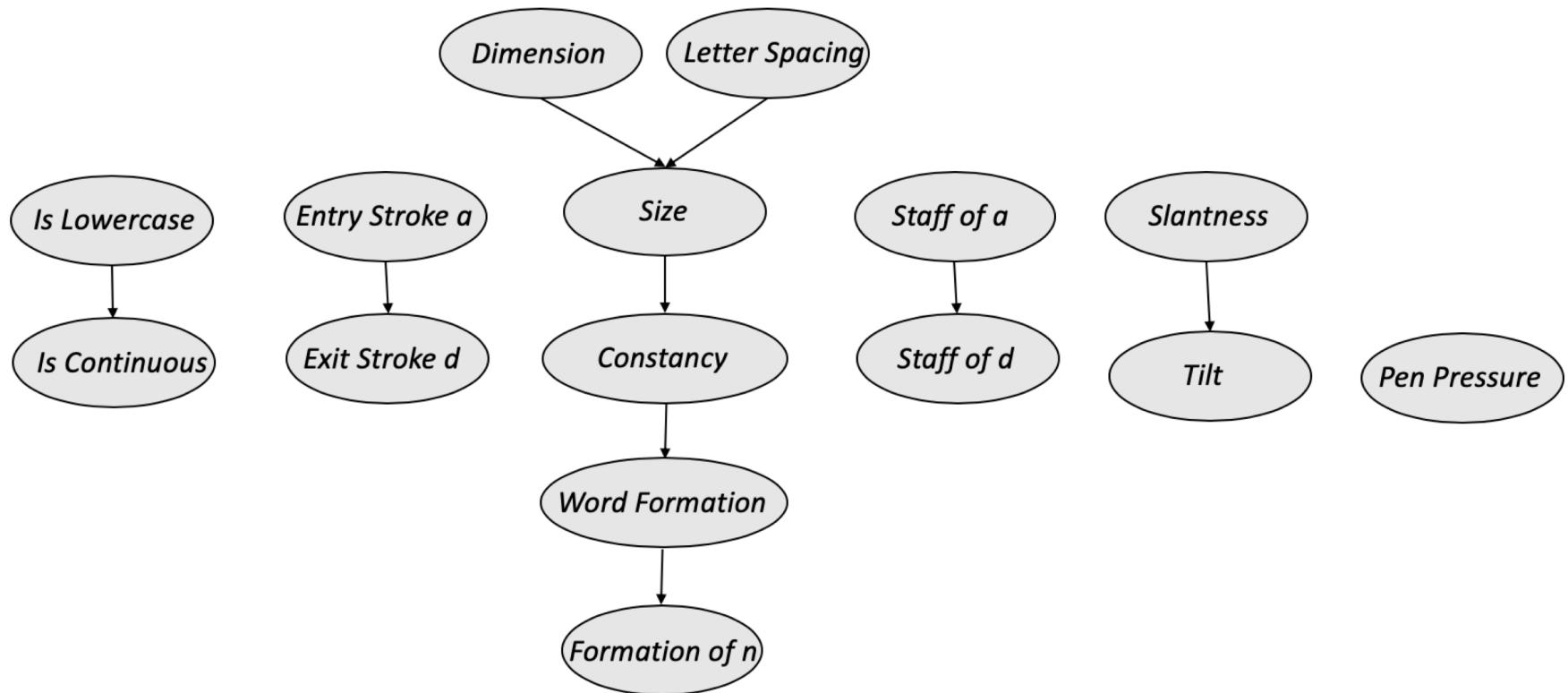


Image 1: Human observed features mapped to image

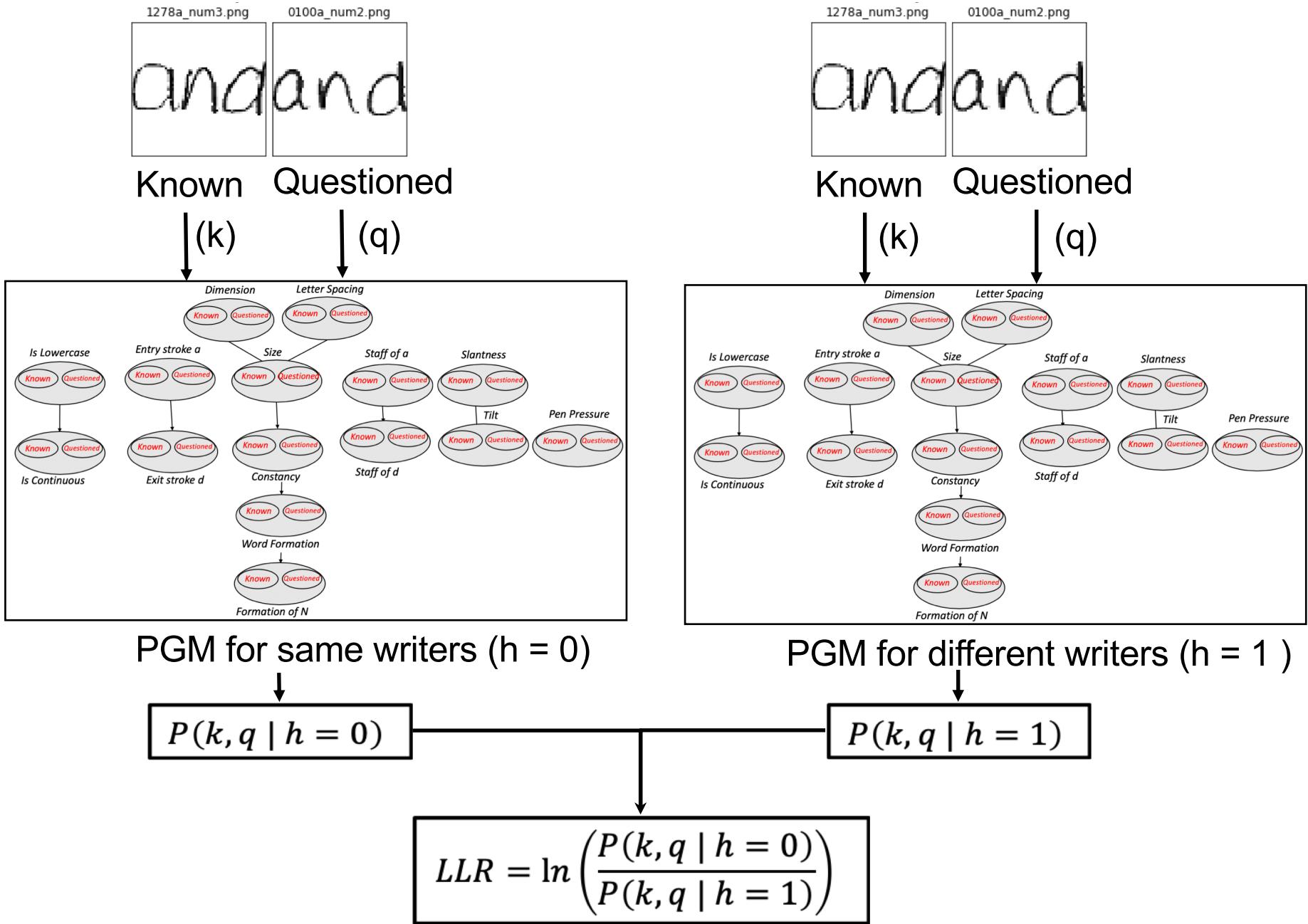
Handwriting comparison – expert features



PGM for human observed ‘and’ dataset



PGM for Writer Verification



Finding LLR

$$LLR = \ln \left(\frac{P(k, q \mid h = 0)}{P(k, q \mid h = 1)} \right)$$

Approach 1

- A simple **Distance based LLR** approach defines scalar distance between known and questioned samples defining loglikelihood ratio as follows:

$$LLR_D = \ln \left(\frac{P(d(k, q) | h = 0)}{P(d(k, q) | h = 1)} \right)$$

Disadvantage: Distribution of N variables is mapped into a scalar thereby incurring severe loss of information (many pairs of k and q can have the same distance).

Approach 2

· **Rarity with Distance based LLR** approach uses strength of the features between the known and the questioned samples defining loglikelihood ratio as follows:

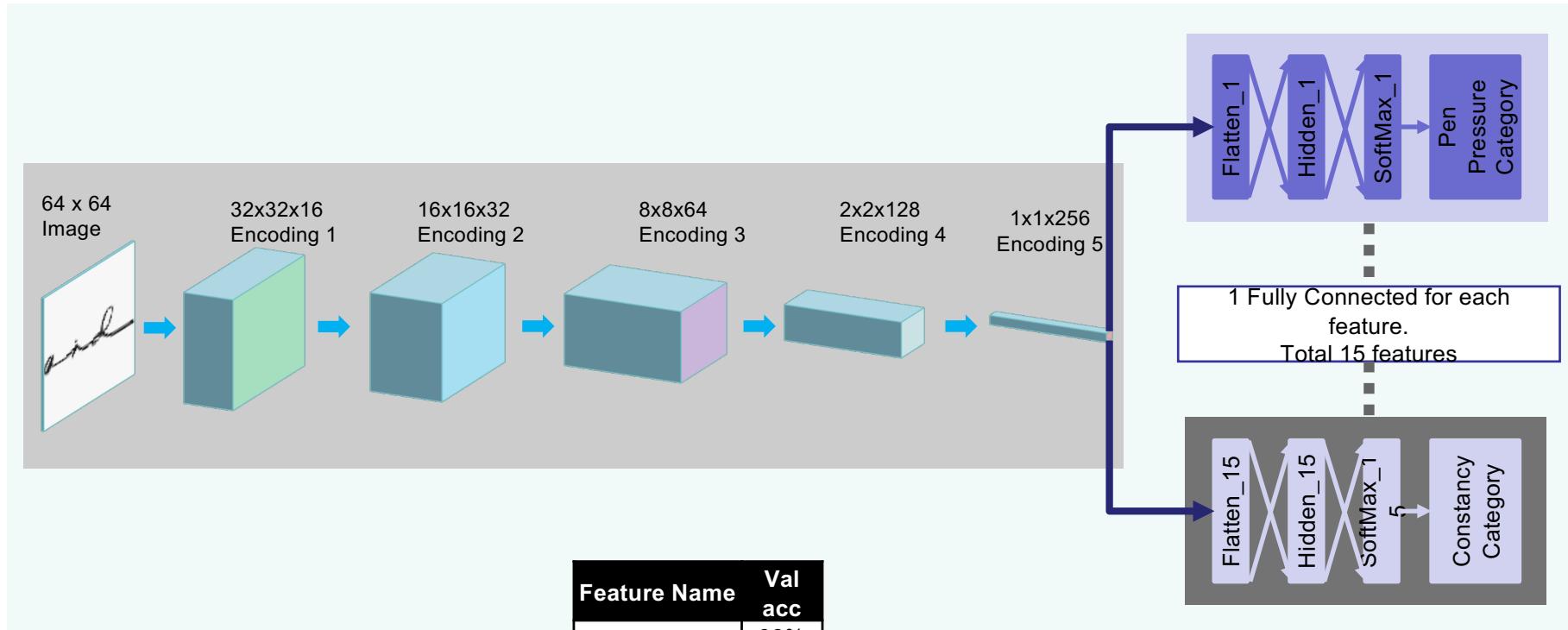
follows:

$$LLR_{RD} = \ln \left(\frac{\frac{1}{\sqrt{2\sigma}} e^{-\frac{(k-q)^2}{4\sigma^2}}}{\frac{1}{\tau} e^{-\frac{(m-\mu)^2}{4\tau^2}}} \right)$$

Distance Rarity

Disadvantage: Not always computationally tractable

Feature extraction architecture/performance



Feature Name	Val acc
pen_pressure	98%
letter_spacing	78%
size	92%
dimension	89%
is_lowercase	99%
is_continuous	94%
slantness	72%
tilt	98%
entry_stroke_a	97%
staff_of_a	84%
formation_n	95%
staff_of_d	86%
exit_stroke_d	66%
wordFormation	89%
constancy	84%

Code: https://github.com/mshaikh2/HDL_Forensics

Comparison performance

Performance based on Overall Similarity

Explainability:

- Compare corresponding 'SoftMax_x' ($x : \{1 \text{ to } 15\}$) of two images using cosine similarity to obtain feature wise explanation

Metric	Seen	Unseen	Shuffled
Intra Writer Accuracy	88%	71%	81%
Inter Writer Accuracy	96%	93%	95%
Average Accuracy	96%	94%	95%

Explanation of Comparison

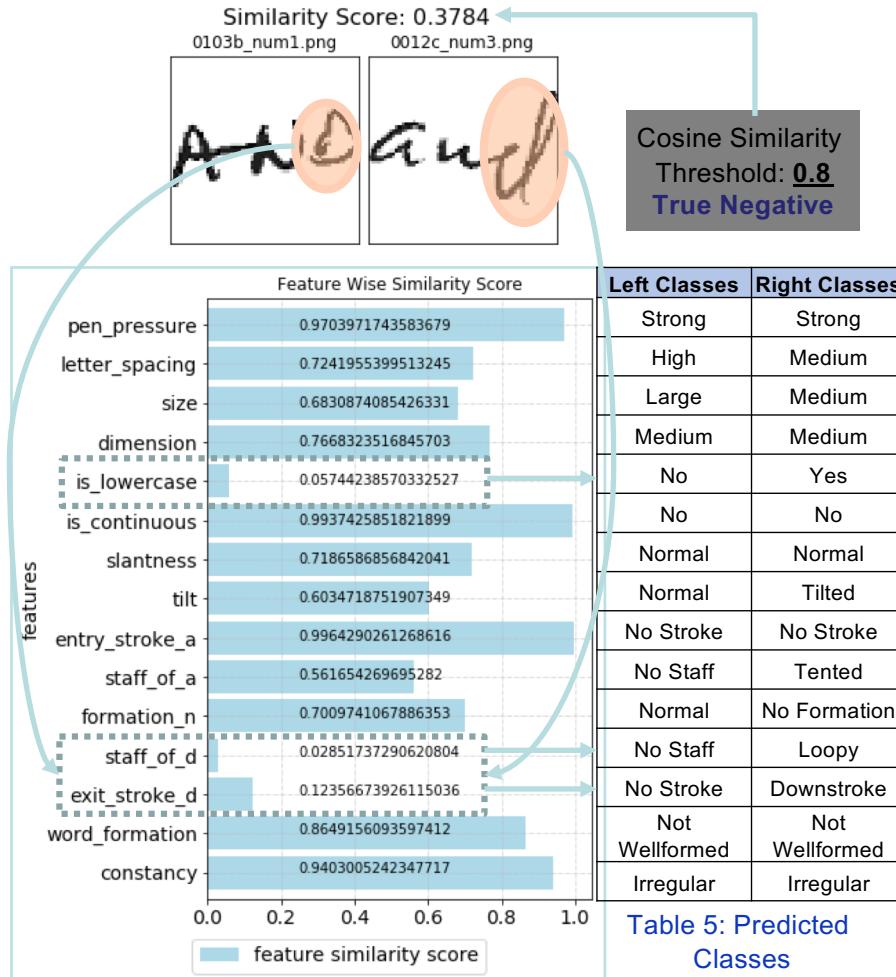


Image 6: Graph of Similarity Score of Left & Right Image Features
Different Writer Samples

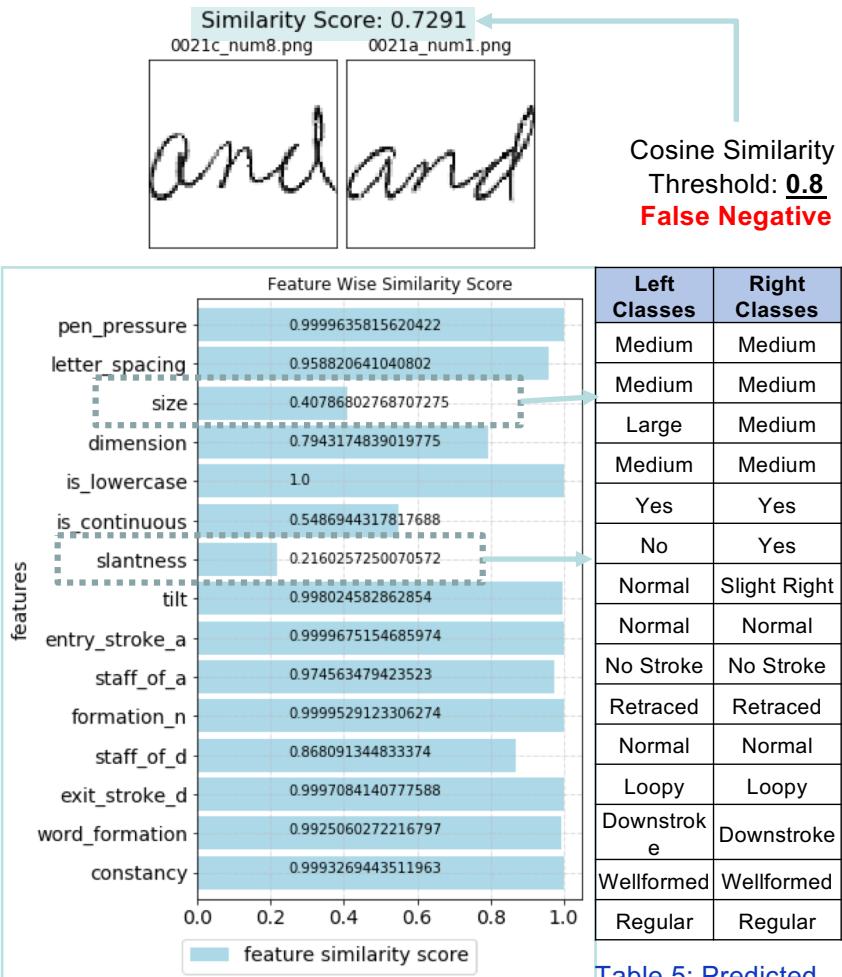
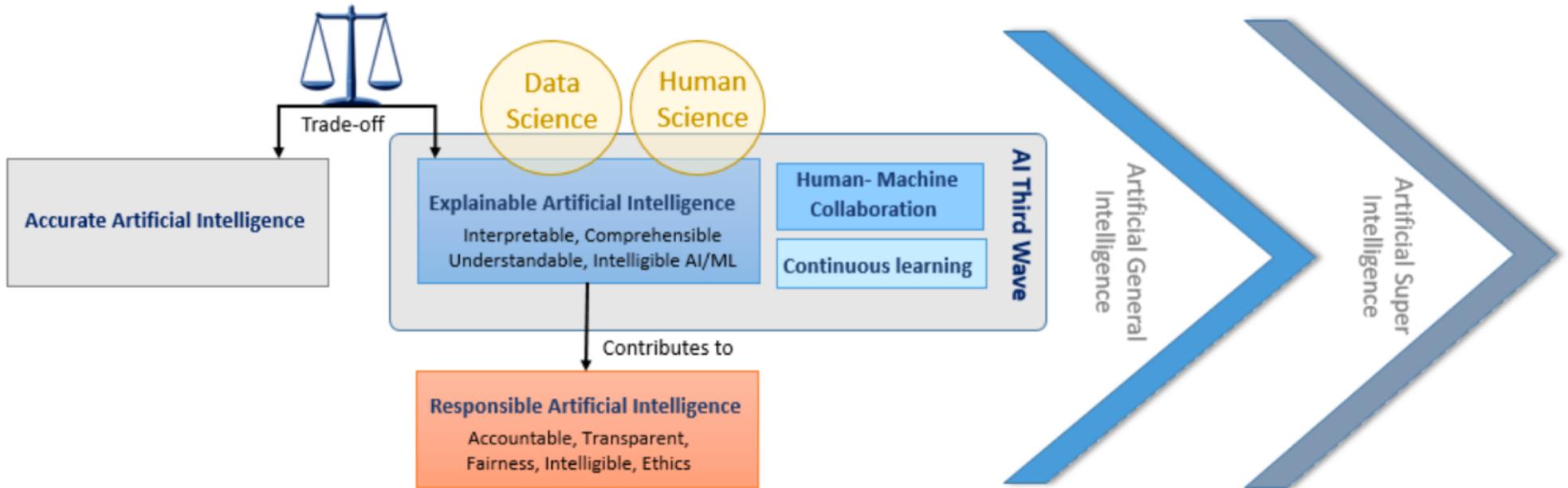


Image 7: Graph of Similarity Score of Left & Right Image Features
Same Writer Samples

Report Card of Comparing two Samples

Future of XAI



Source: <https://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=8466590>

Summary and Conclusion

- Deep Learning models are often blackboxes
- Ante-hoc XAI models are new learning models with an explanatory interface
- Post-hoc models retain high performance
- Several new ideas developing: Bayesian Teaching, Causal learning of deep neural nets
- Probabilistic AI is useful for deep learning
- XAI may be an inflection point in AI