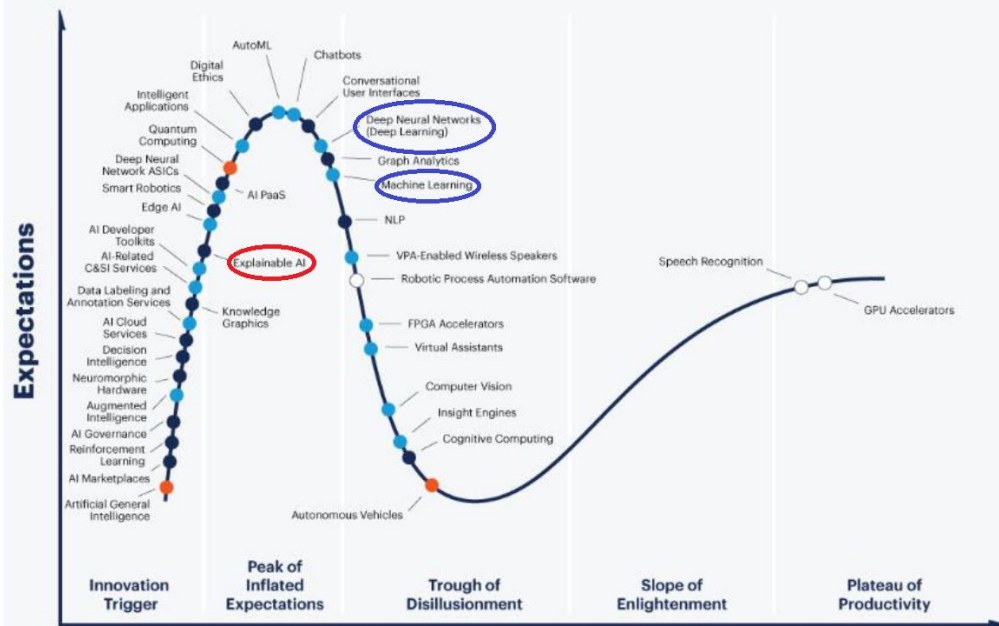


# Объясняемый искусственный интеллект (ХАИ)

Н.Ю. Золотых  
ННГУ, ИТММ

2 октября 2020

# Gartner Hype Cycle for Artificial Intelligence, 2019



Plateau will be reached:

○ less than 2 years

● 2 to 5 years

● 5 to 10 years

● more than 10 years

● obsolete before plateau

As of July 2019

[gartner.com/SmarterWithGartner](https://gartner.com/SmarterWithGartner)

Source: Gartner  
© 2019 Gartner, Inc. and/or its affiliates. All rights reserved.

**Gartner**

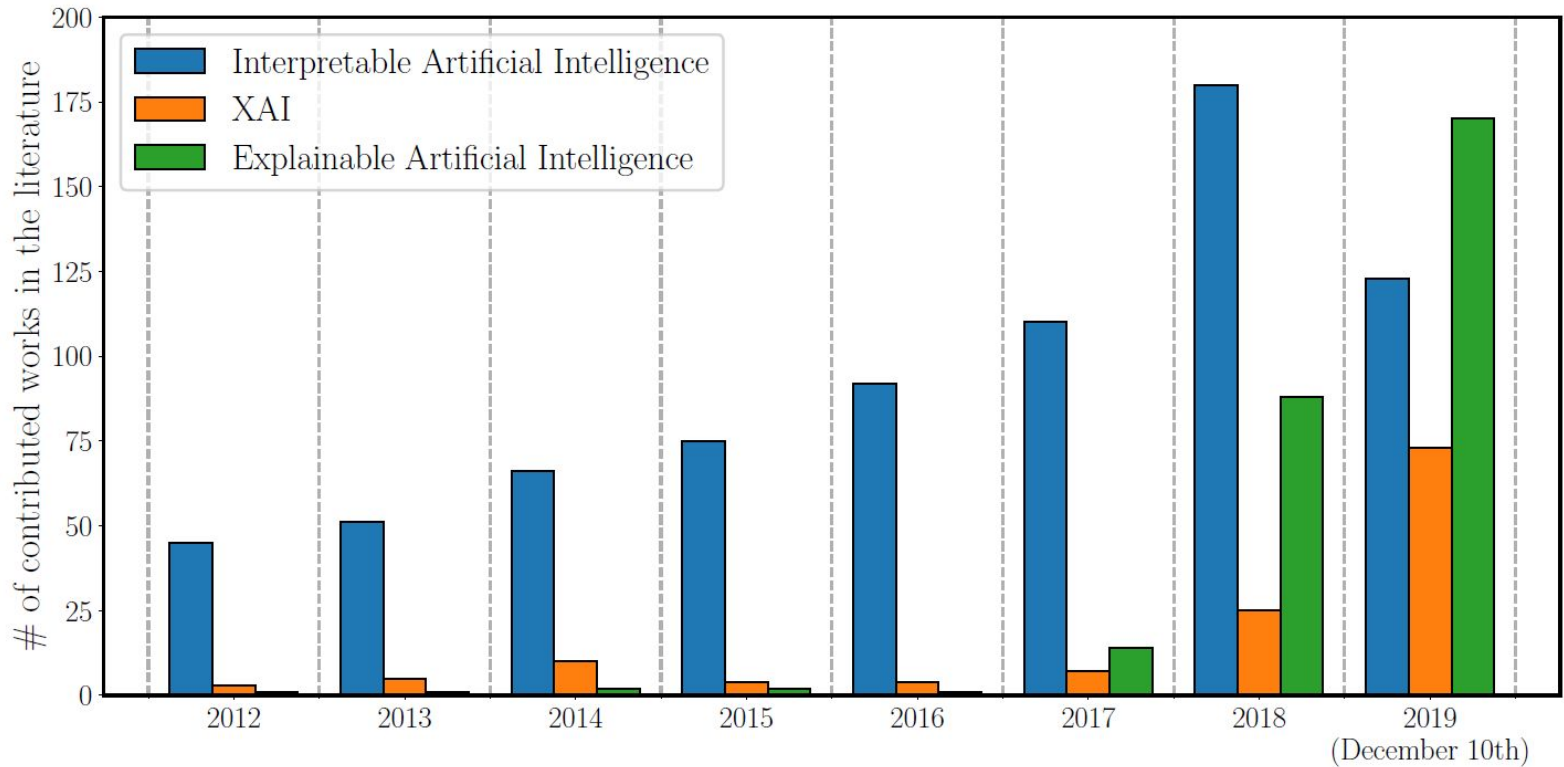


Figure 1: Evolution of the number of total publications whose title, abstract and/or keywords refer to the field of XAI during the last years. Data retrieved from Scopus<sup>®</sup> (December 10th, 2019) by using the search terms indicated in the legend when querying this database. It is interesting to note the latent need for interpretable AI models over time (which conforms to intuition, as interpretability is a requirement in many scenarios), yet it has not been until 2017 when the interest in techniques to explain AI models has permeated throughout the research community.

- *D. Gunning*, Explainable artificial intelligence (xAI), Tech. rep., Defense Advanced Research Projects Agency (DARPA) (2017)
- *A. B. Arrieta*, et al. Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion* 58 (2020): 82-115.
- *E. Tjoa, C. Guan*, A survey on explainable artificial intelligence (XAI): Towards medical XAI (2019). arXiv:1907.07374.
- *L. H. Gilpin, D. Bau, B. Z. Yuan, A. Bajwa, M. Specter, L. Kagal*, Explaining Explanations: An Overview of Interpretability of Machine Learning (2018). arXiv:1806.00069
- *F. K. Došilović, M. Brčić, N. Hlupić*, Explainable artificial intelligence: A survey, in: 41st International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO), 2018, pp. 210–215.
- *A. Adadi, M. Berrada*, Peeking inside the black-box: A survey on explainable artificial intelligence (XAI), *IEEE Access* 6 (2018) 52138–52160.
- *O. Biran, C. Cotton*, Explanation and justification in machine learning: A survey, in: *IJCAI-17 workshop on explainable AI (XAI)*, Vol. 8, 2017, p. 1.
- *S. T. Shane, T. Mueller, R. R. Hoffman, W. Clancey, G. Klein*, Explanation in Human-AI Systems: A Literature Meta-Review Synopsis of Key Ideas and Publications and Bibliography for Explainable AI, Tech. rep., Defense Advanced Research Projects Agency (DARPA) XAI Program (2019).
- *R. Guidotti, A. Monreale, S. Ruggieri, F. Turini, F. Giannotti, D. Pedreschi*, A survey of methods for explaining black box models, *ACM Computing Surveys* 51 (5) (2018) 93:1–93:42.

# Interpretability & Explainability

*Интерпретация* - это отображение абстрактного понятия в понятную для человека область [Montavon, Samek, Muller, 2018]

*Объяснение* - это набор характеристик интерпретируемой области, которые привели к принятию данного решения в конкретном случае [Montavon, Samek, Muller, 2018]

Comprehensibility  $\approx$  interpretability  $\approx$  model-centric explainability

Transparency  $\approx$  explainability  $\approx$  subject-centric explainability

# Machine Learning vs Data Mining

Говорят, что компьютерная программа *обучается* на основе опыта  $E$  по отношению к некоторому классу задач  $T$  и меры качества  $P$ , если качество решения задач из  $T$ , измеренное на основе  $P$ , улучшается с приобретением опыта  $E$ .

*T.M. Mitchell, 1997*

*Data Mining* – совокупность методов обнаружения в данных ранее неизвестных, нетривиальных, практически полезных и *доступных интерпретации знаний*, необходимых для принятия решений в различных сферах человеческой деятельности.

*Г. Пятецкий-Шапиро, 1989*

# Machine Learning vs Data Mining

ML и DM извлекают закономерности («знания») из данных, но (немного) с разными целями:

- ML – чтобы обучить машину;
- DM – чтобы обучить человека.

Поэтому *в первую очередь\**

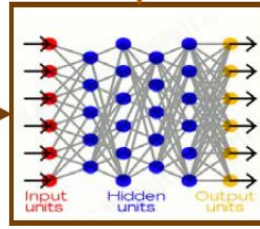
- в ML минимизируют ошибку;
- в DM важна интерпретируемость результата.

XAI должен стереть эту границу

# Today



Training Data



Learned Function

**This is a cat**  
(p = .93)

Output



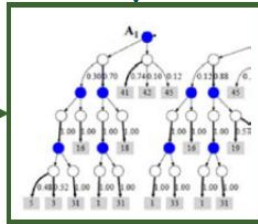
User with a Task

- Why did you do that?
- Why not something else?
- When do you succeed?
- When do you fail?
- When can I trust you?
- How do I correct an error?

# Tomorrow



Training Data



Explainable Model

**This is a cat:**  
• It has fur, whiskers, and claws.  
• It has this feature:



Explanation Interface



User with a Task

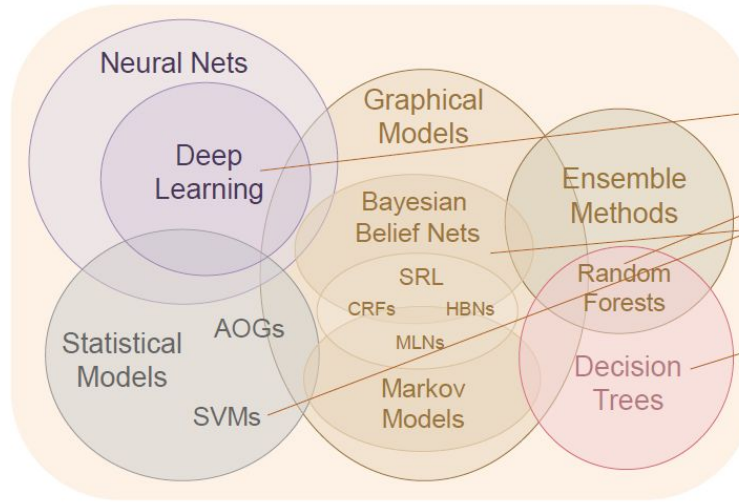
- I understand why
- I understand why not
- I know when you'll succeed
- I know when you'll fail
- I know when to trust you
- I know why you erred



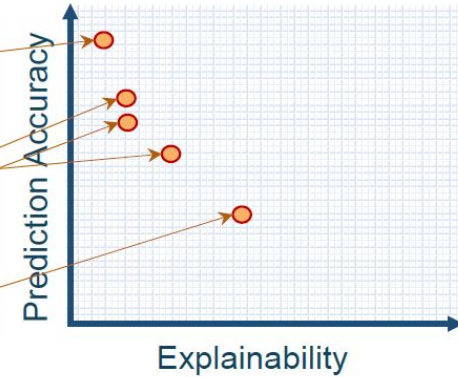
## New Approach

Create a suite of machine learning techniques that produce more explainable models, while maintaining a high level of learning performance

## Learning Techniques (today)



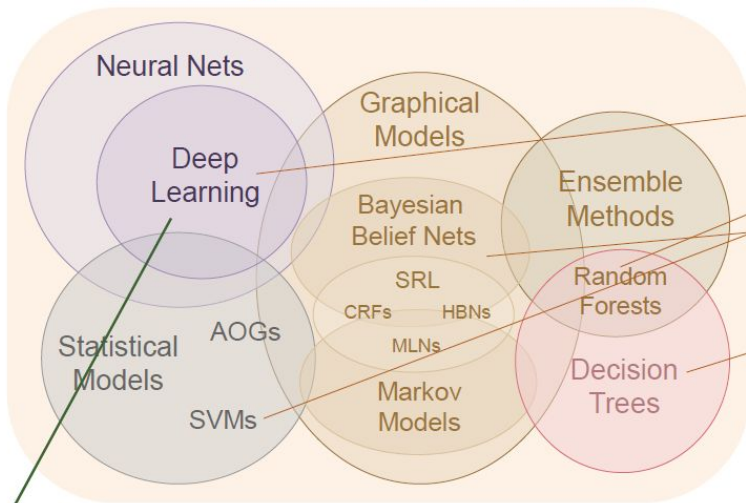
## Explainability (notional)



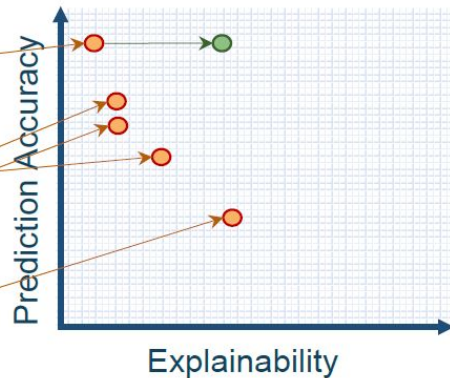
## New Approach

Create a suite of machine learning techniques that produce more explainable models, while maintaining a high level of learning performance

## Learning Techniques (today)



## Explainability (notional)



A diagram of a deep neural network with three layers of nodes. The input layer has red nodes, the hidden layer has blue nodes, and the output layer has yellow nodes. Below the network, the words 'Whiskers' and 'Claws' are shown in ovals, with arrows pointing to the corresponding nodes in the network. The diagram is labeled 'Deep Explanation' and 'Modified deep learning techniques to learn explainable features'.

**Deep Explanation**  
Modified deep learning techniques to learn explainable features

# Deep Explanation

- Извлечение знаний путем комплексного упрощения сети (А.Н. Горбань, 1990, В.Г. Царегородцев, 1998)
- Контрастирование и логически прозрачные нейронные сети (А.Н. Горбань, 1990, Д.И.Еремин, Е.М. Миркес, 1994)
- Визуализация слоев сверточной сети с помощью деконволюции (Zeiler, Fergus, 2014)
- Интерпретация узлов сети как семантических понятий, например, как при идентификации элементов на изображении или регистрации событий (Yu, Liu, Cheng, Divakaran, Sawhney, 2012; Gan, Wang, Ян, Юнг, Хауптманн, 2015)
- Использование методов создания подписей к изображениям (LeCun, Bengio, Hinton, 2015) для генерации объяснений (Hendricks et al., 2016)
- ...

## Multimedia Event Recounting

Ranked Videos

Expanded View

Primary Evidence

Bride walking with a man with people watching

Bride and groom with officiant

Bride and groom put rings on hands

Evidence Composition

Scene: Indoors :: Church

Walking together

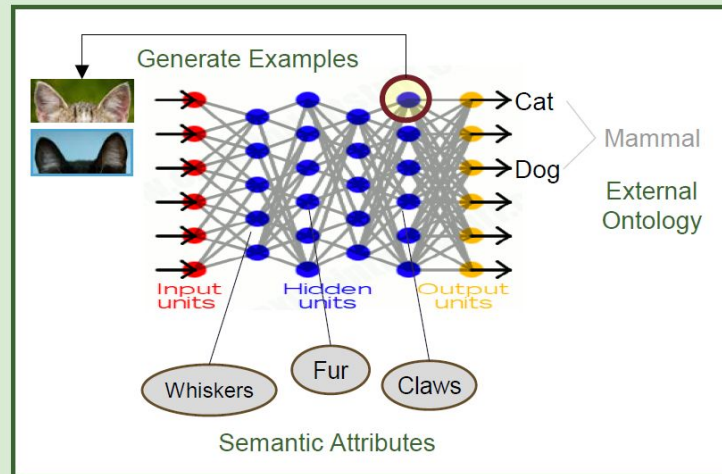
A number of faces detected  
→ infer Group of people

person: bride

person: bride

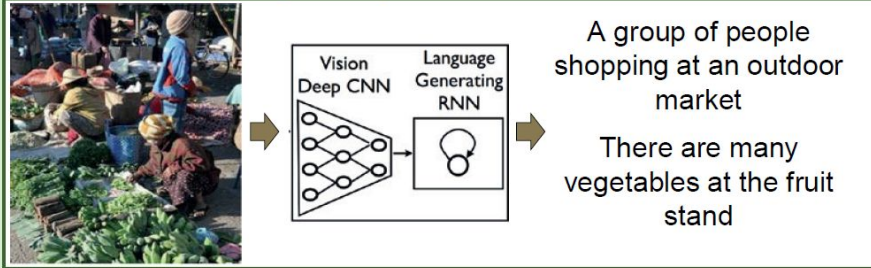
- This illustrates an example of event recounting.
- The system classified this video as a wedding.
- The frames above show its evidence for the wedding classification

## Learning Semantic Associations



- Train the net to associate semantic attributes with hidden layer nodes
- Train the net to associate labelled nodes with known ontologies
- Generate examples of prominent but unlabeled nodes to discover semantic labels
- Generate clusters of examples from prominent nodes
- Identify the best architectures, parameters, and training sequences to learn the most interpretable models

## Generating Image Captions

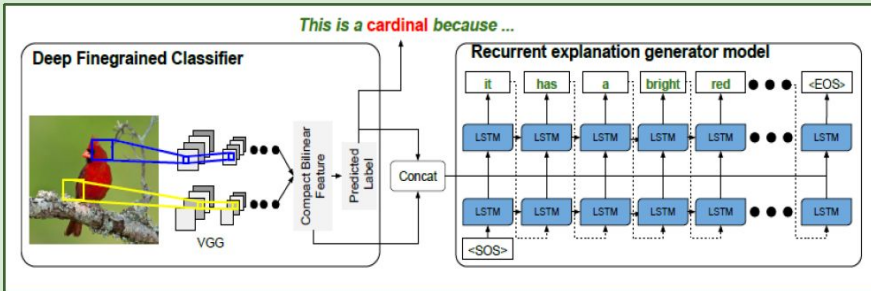


- A CNN is trained to recognize objects in images
- A language generating RNN is trained to translate features of the CNN into words and captions.

## Example Explanations



## Generating Visual Explanations



Researchers at UC Berkeley have recently extended this idea to generate explanations of bird classifications. The system learns to:

- Classify bird species with 85% accuracy
- Associate *image descriptions* (discriminative features of the image) with *class definitions* (image-independent discriminative features of the class)

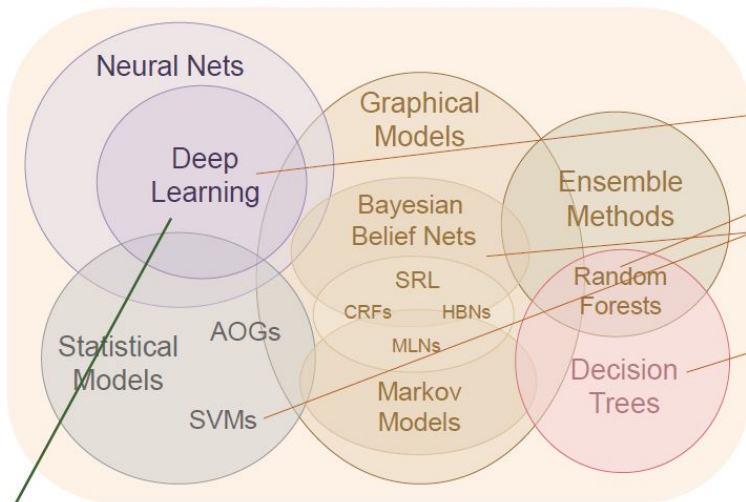
## Limitations

- Limited (indirect at best) explanation of internal logic
- Limited utility for understanding classification errors

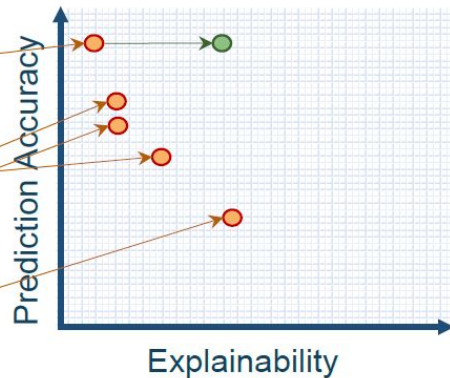
## New Approach

Create a suite of machine learning techniques that produce more explainable models, while maintaining a high level of learning performance

## Learning Techniques (today)



## Explainability (notional)



A diagram of a deep neural network with three layers of nodes. The input layer has nodes labeled 'Input Layer', the hidden layer has nodes labeled 'Hidden Layer', and the output layer has nodes labeled 'Output Layer'. Below the network, two input features are shown: 'Whiskers' and 'Claws'. The network is used to illustrate the concept of 'Deep Explanation'.

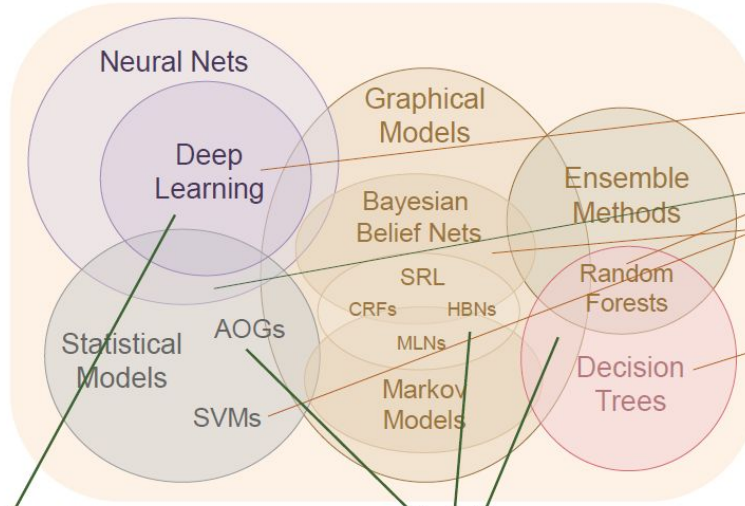
### Deep Explanation

Modified deep learning techniques to learn explainable features

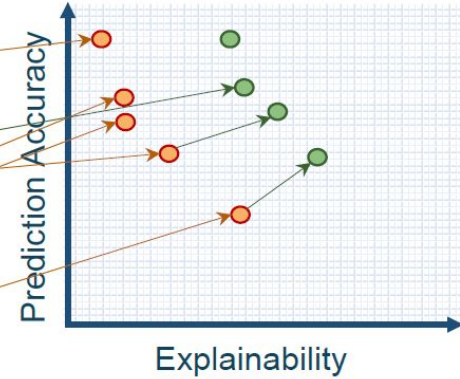
# New Approach

Create a suite of machine learning techniques that produce more explainable models, while maintaining a high level of learning performance

# Learning Techniques (today)



# Explainability (notional)



A diagram of a deep neural network with three hidden layers. The input layer has nodes for 'Whiskers' and 'Claws'. The hidden layers are labeled 'Hidden Layer 1' and 'Hidden Layer 2'. The output layer has nodes for 'Cat' and 'Not Cat'. The network is shown with various weights and connections between nodes.

## Deep Explanation

Modified deep learning techniques to learn explainable features

A decision tree diagram with a root node  $A_1$ . The tree branches into several nodes, each with numerical values and labels. The nodes are labeled with 'Yes' and 'No' outcomes. The tree is shown with various weights and connections between nodes.

## Interpretable Models

Techniques to learn more structured, interpretable, causal models

# Interpretable Models

Структурированные, интерпретируемые и причинно-следственные модели:

- Bayesian Rule Lists (Letham, Rudin, McCormick, Madigan, 2015)
- генеративные модели, такие, как Bayesian Program Learning (Lake, Salakhutdinov, Tenenbaum, 2015)
- использование стохастических грамматик (Brendel Todorovic, 2011; Park, Nie, Zhu 2016)
- модели причинно-следственных связей (Maier, Taylor, Oktay, Jensen, 2010)
-



Training Data  
1623 Characters



Bayesian  
Program  
Learning

```

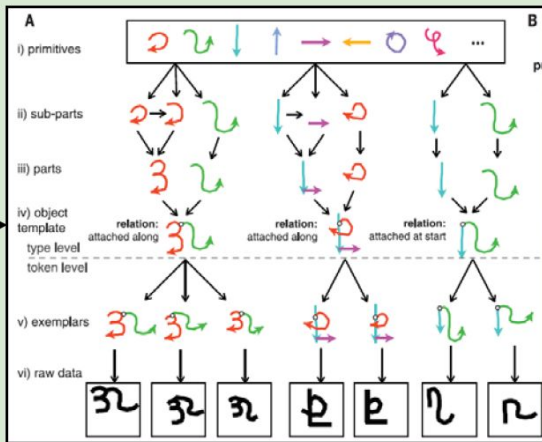
num_strokes ~ Poisson(2)
for i = 1 to num_strokes:
  num_substrokes_prior[i] ~ Discrete([0,1,1,1])
  num_substrokes[i] ~ Poisson(num_substrokes_prior[i])
  for j = 1 to num_substrokes[i]:
    substrokes[i][j] ~ substroke_transition_prob(i)->
      relation[i] ~ relation_prob(substrokes[i])

for i = 1 to num_strokes:
  noised_substrokes[i][:] = stroke_noise(substrokes[i][:])
  stroke_start_position[i] = start_distribution(relation,
  trajectory[i]) = draw_trajectory(stroke_start_position[i],
  AffineTransform = transform_distribution
  image = render(AffineTransform(trajectory))
  
```

Seed Model

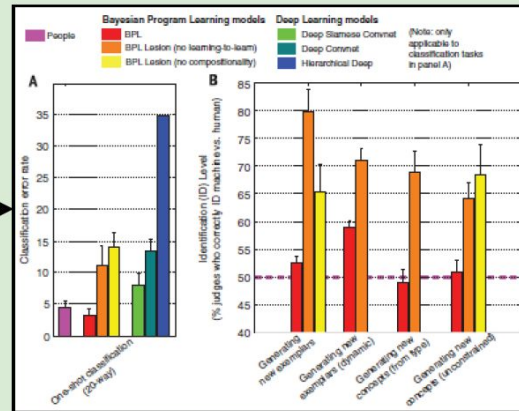
A simple Probabilistic Program that describes the parameters of character generation

# Concept Learning Through Probabilistic Program Induction



Generative Model

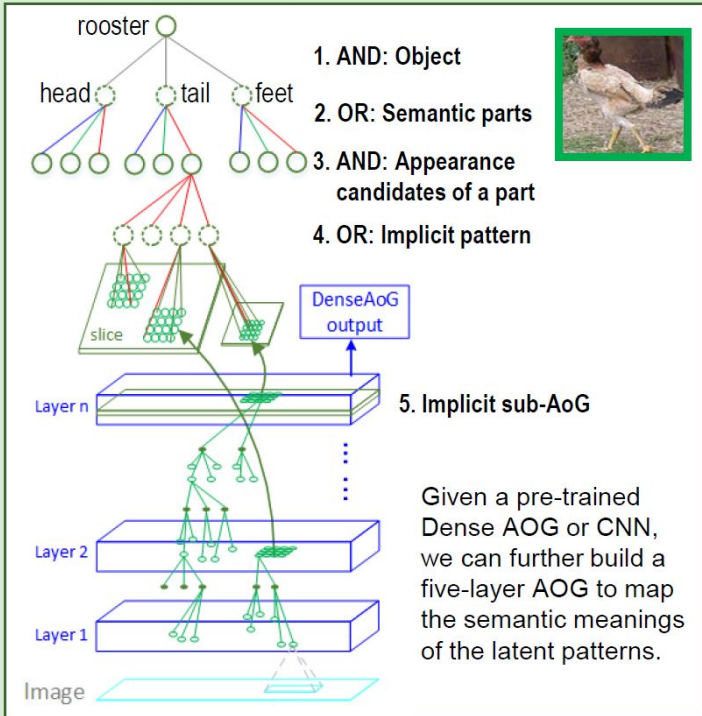
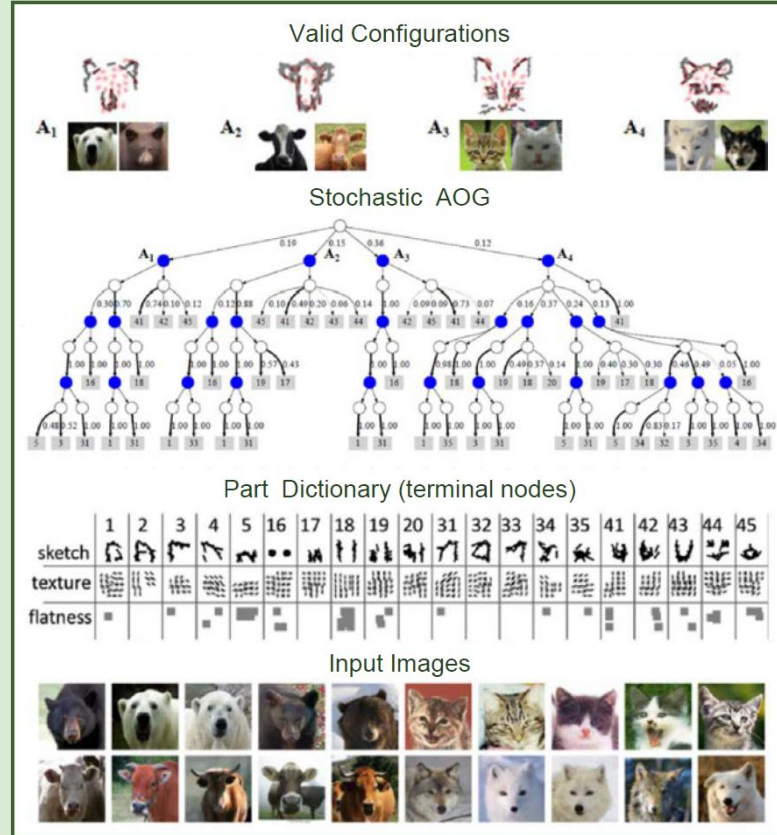
Recognizes characters by generating an explanation of how a new test character might be created (i.e., the most probable sequence of strokes that would create that character)



Performance

This model matches human performance and out performs deep learning

# Stochastic And-Or-Graphs (AOG)

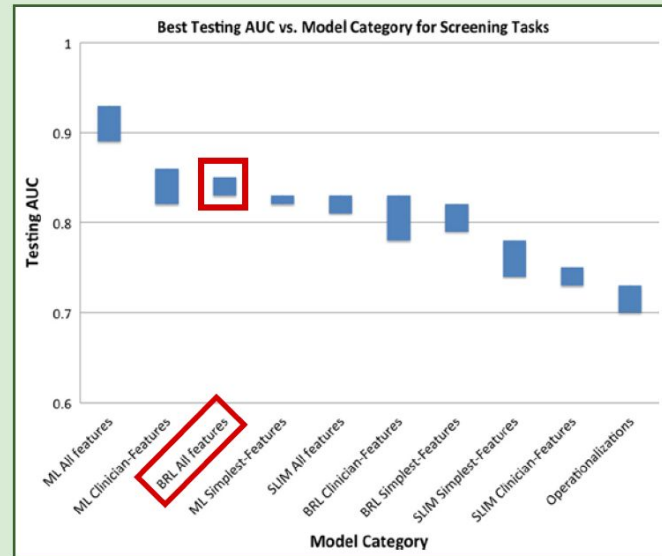
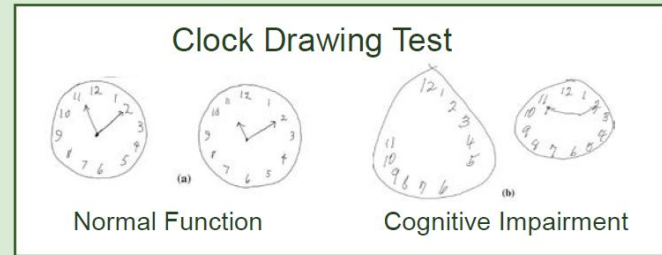


$$L(\theta) = \frac{1}{M} \sum_{m=1}^M \underbrace{\log P(I_m, \theta)}_{\text{generative}} + \underbrace{L(pg_m^*, \hat{p}g_m)}_{\text{discriminative}}$$

# Bayesian Rule Lists (BRL)

- **if** hemiplegia and age > 60
  - **then** stroke risk 58.9% (53.8%–63.8%)
- **else if** cerebrovascular disorder
  - **then** stroke risk 47.8% (44.8%–50.7%)
- **else if** transient ischaemic attack
  - **then** stroke risk 23.8% (19.5%–28.4%)
- **else if** occlusion and stenosis of carotid artery without infarction
  - **then** stroke risk 15.8% (12.2%–19.6%)
- **else if** altered state of consciousness and age > 60
  - **then** stroke risk 16.0% (12.2%–20.2%)
- **else if** age ≤ 70
  - **then** stroke risk 4.6% (3.9%–5.4%)
- **else** stroke risk 8.7% (7.9%–9.6%)

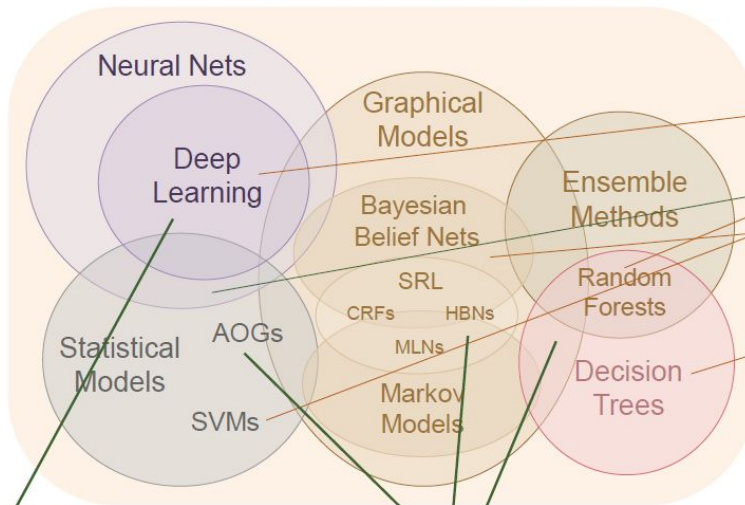
- BRLs are decision lists—a series of if-then statements
- BRLs discretize a high-dimensional, multivariate feature space into a series of simple, readily interpretable decision statements.
- Experiments show that BRLs have predictive accuracy on par with the current top ML algorithms (approx. 85–90% as effective) but with models that are much more interpretable



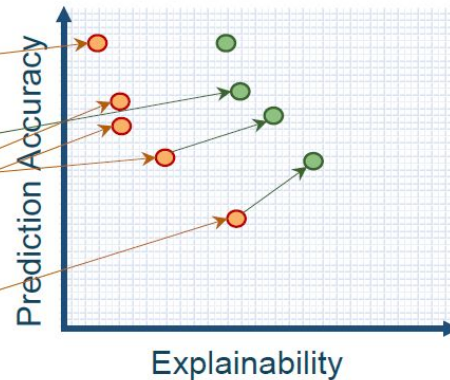
## New Approach

Create a suite of machine learning techniques that produce more explainable models, while maintaining a high level of learning performance

## Learning Techniques (today)



## Explainability (notional)



A diagram of a deep neural network with three hidden layers. The input layer has nodes labeled 'Whiskers' and 'Claws'. The hidden layers are labeled 'Hidden Layer 1' and 'Hidden Layer 2'. The output layer has nodes labeled 'Fur' and 'Claws'. The network is used to illustrate 'Deep Explanation'.

### Deep Explanation

Modified deep learning techniques to learn explainable features

A diagram of a decision tree model. The root node is labeled  $A_1$ . The tree branches out into several nodes, each containing numerical values and labels like 'Yes' and 'No'. The tree is used to illustrate 'Interpretable Models'.

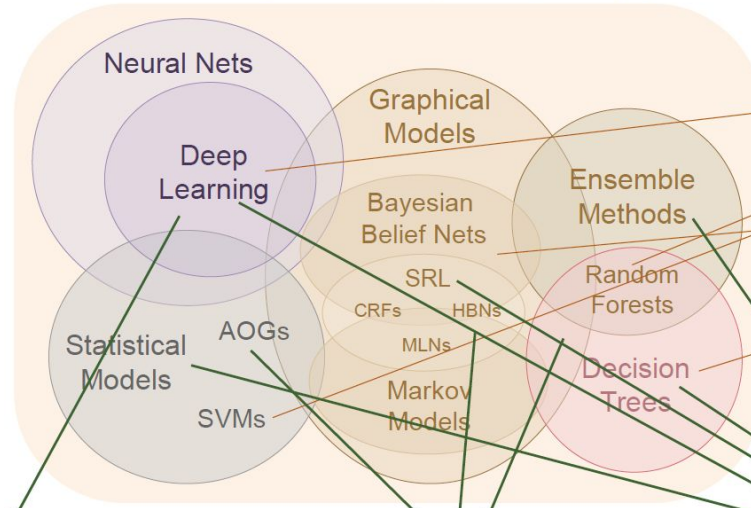
### Interpretable Models

Techniques to learn more structured, interpretable, causal models

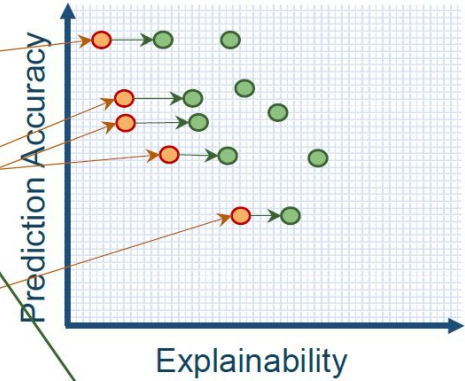
## New Approach

Create a suite of machine learning techniques that produce more explainable models, while maintaining a high level of learning performance

## Learning Techniques (today)



## Explainability (notional)



**Deep Explanation**  
Modified deep learning techniques to learn explainable features

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

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

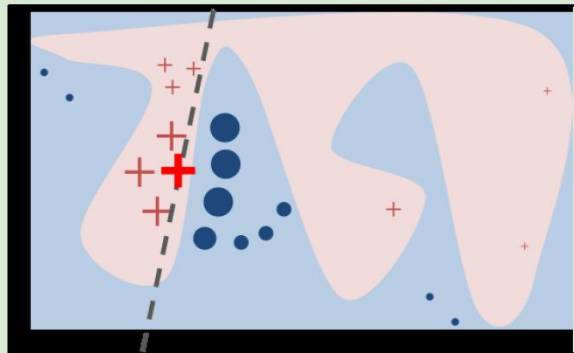
# Model Induction

Эксперименты с моделями как с черными ящиками (model-agnostic explanation system).

- исследование откликов на разные входы (например, Ribeiro, Singh, & Guestrin 2016)
- использование абдукции, рассуждений и story generation для “рационализации” правдоподобных объяснений поведения системы
- ...

# Local Interpretable Model-agnostic Explanations (LIME)

## Black-box Induction



The black-box model's complex decision function  $f$  (unknown to LIME) is represented by the blue/pink background. The bright bold red cross is the instance being explained. LIME samples instances, gets predictions using  $f$ , and weighs them by the proximity to the instance being explained (represented here by size). The dashed line is the learned explanation that is locally (but not globally) faithful. .

## Example Explanation



(a) Original Image

Electric Guitar  
 $p = 0.32$



(b) Explaining *Electric guitar*

Acoustic Guitar  
 $p = 0.24$



(c) Explaining *Acoustic guitar*

- **LIME** is an algorithm that can explain the predictions of any classifier in a faithful way, by approximating it locally with an interpretable model.
- **SP-LIME** is a method that selects a set of representative instances with explanations as a way to characterize the entire model.

Спасибо за внимание!