# Essence of Machine Learning (and Deep Learning)

Hoa M. Le

Data Science Lab, HUST

[hoamle.github.io](hoamle.github.io)

# Examples

- https://www.youtube.com/watch?v=BmkA1ZsG2P4

- http://www.r2d3.us/visual-intro-to-machine-learning-part-1/

# Machine Learning is about ...

… a computer program (machine) *learns* to do a task (problem) from experience (data)

- *learning* ≜ improved *performance* with more experience

<div align="right">

*- Tom Mitchell*

</div>

⇑

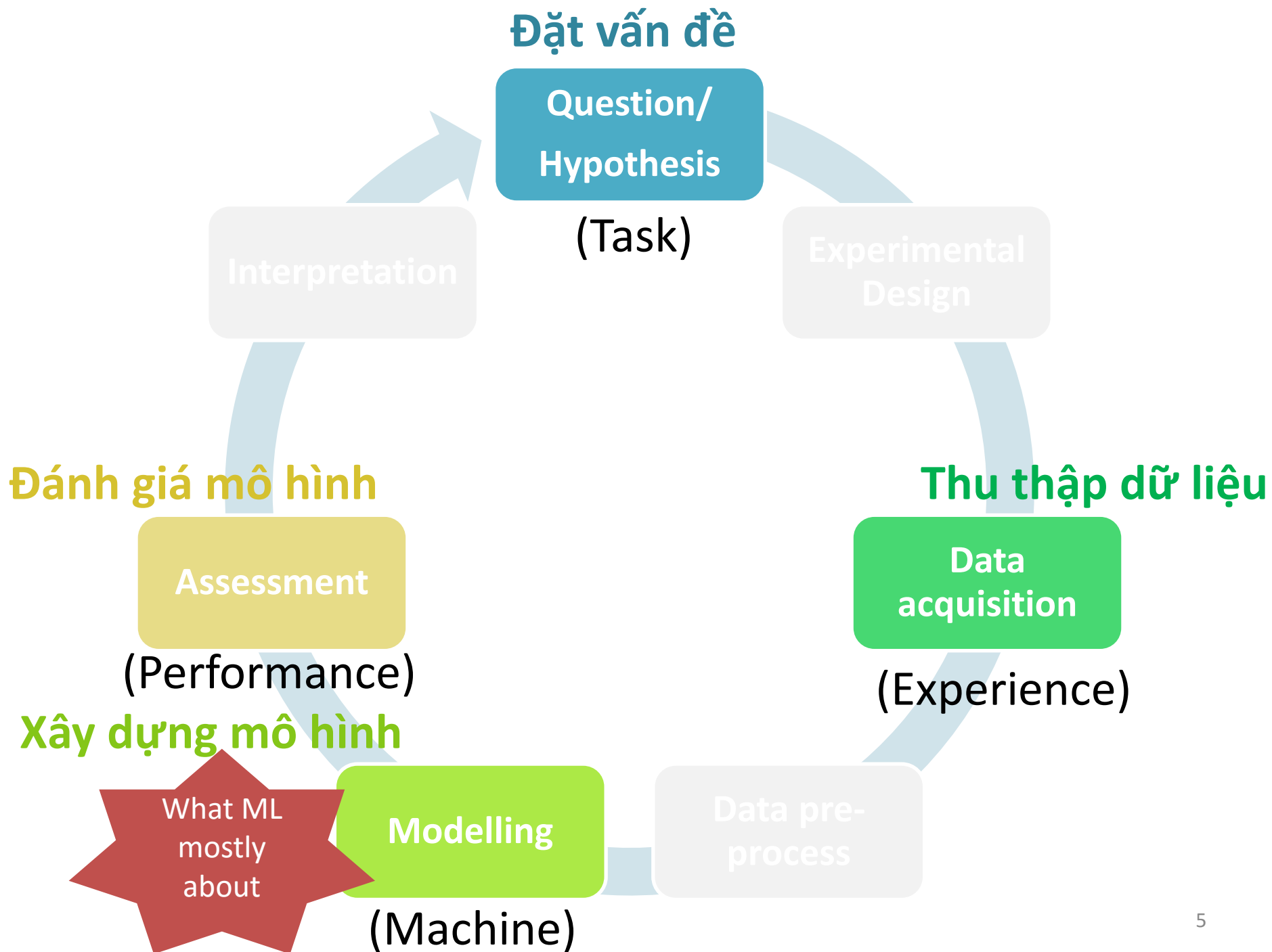**predictive modelling** with **sample data**
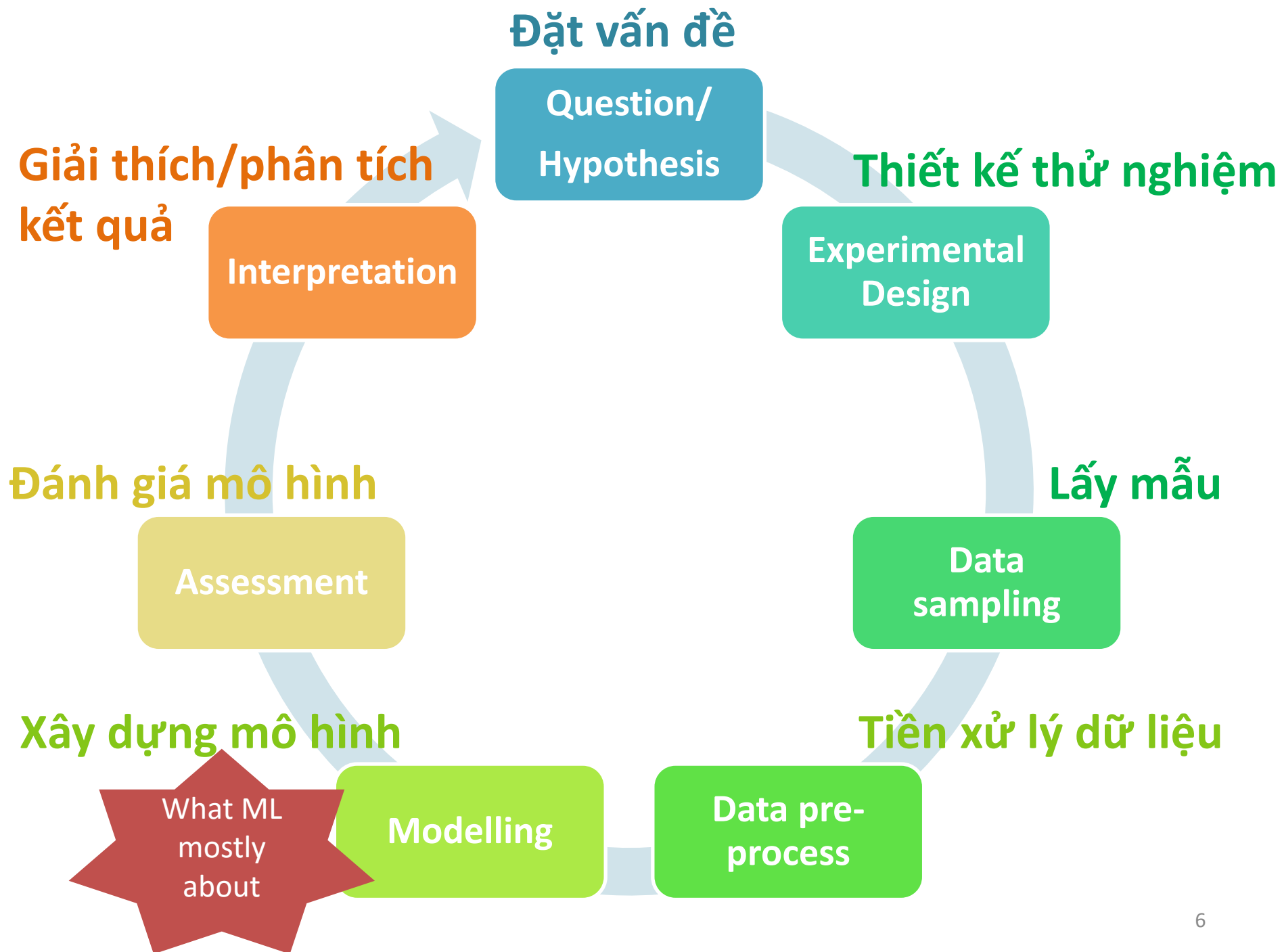
⇑

"heurestics" & statistical modelling

note 1: "heurestic" as in "intuitive, but not (yet!) rigorously proven by mathematical tools at some extend"

*note 2: predictive modelling can also be in the form of rule-based systems, models in physics, etc*

# BUILD
# A MACHINE LEARNING SOLUTION

the Pipeline

# Đặt vấn đề

**Question/ Hypothesis**

**Q.a. What are** there in an **abitrary photo?**
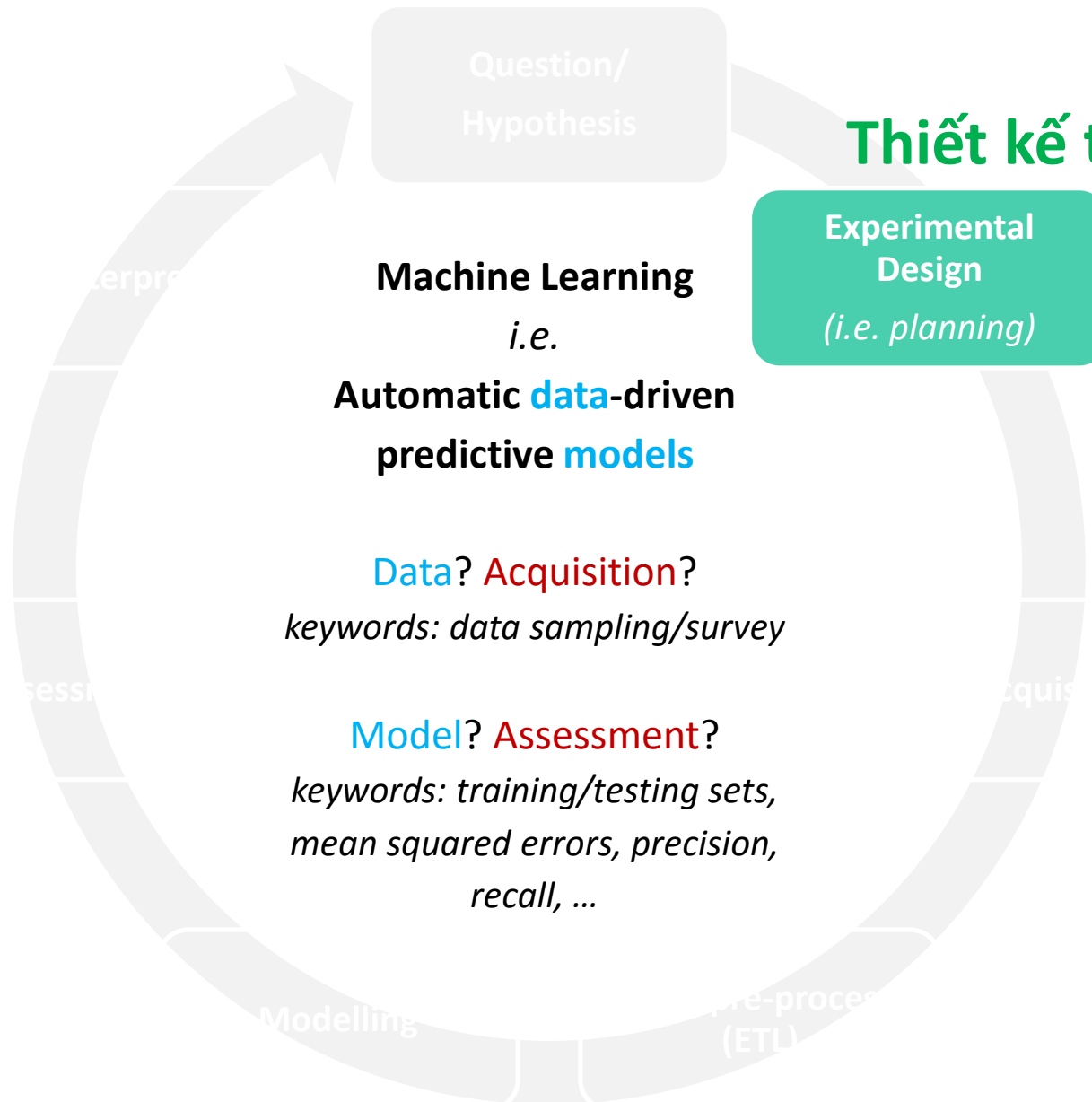**Q.b. What is** there in an **abitrary photo?**
**Q.c. Is there** any puppy an **abitrary photo?**

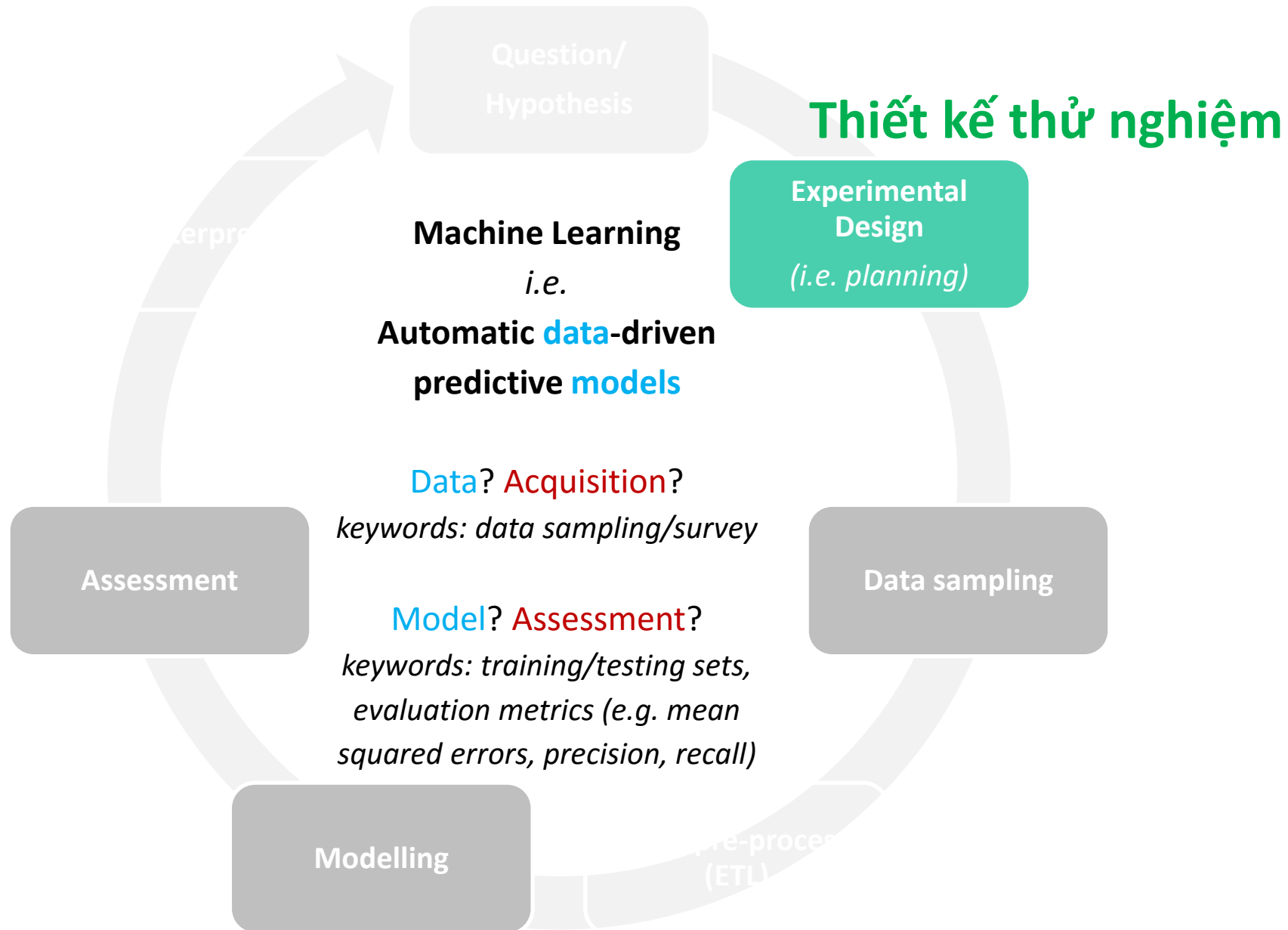

cat
flower
dog
jet
ground
grass
…

*Other questions:*
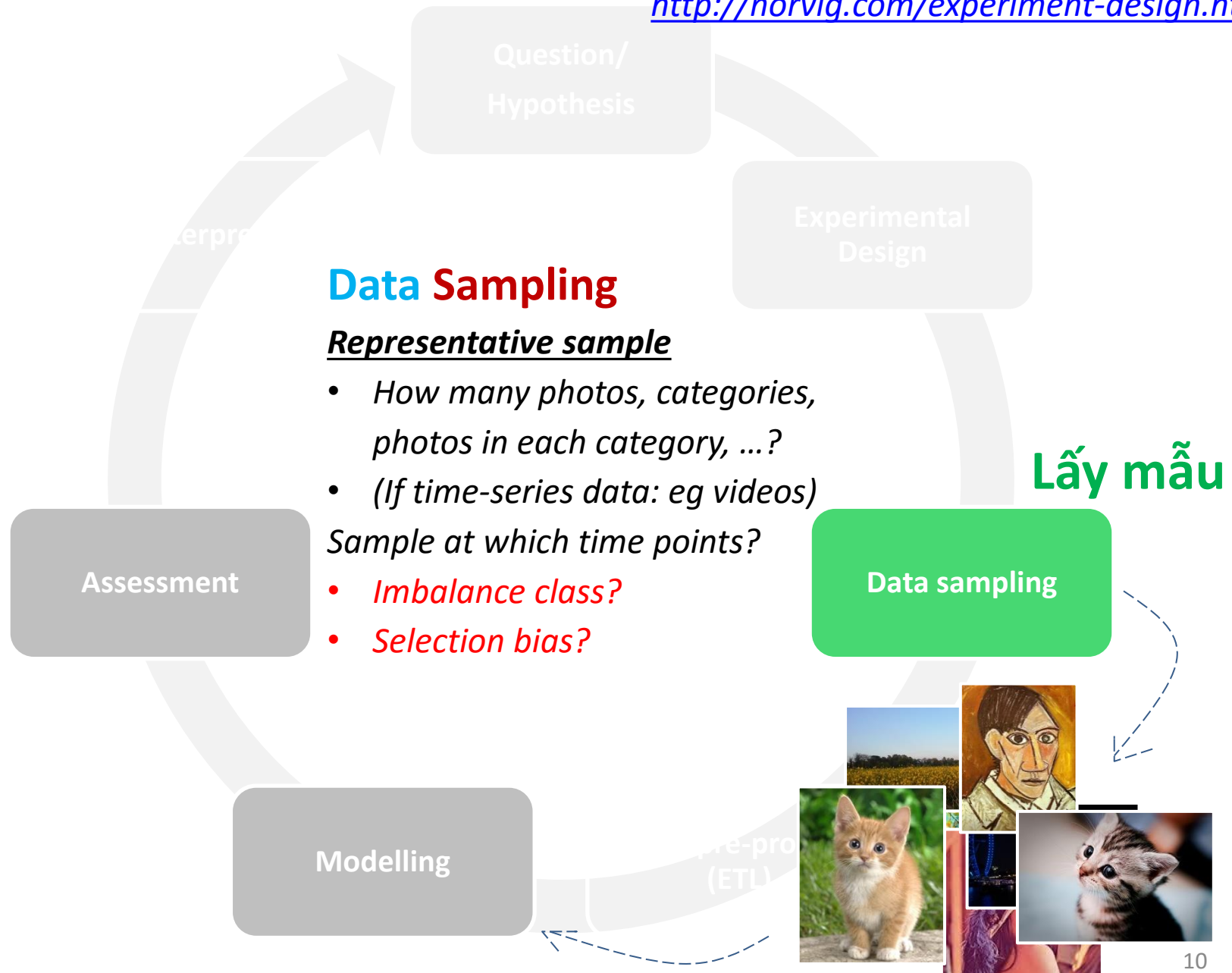- *Where are the puppies in a photo?*
- *How confident can I assure that there is a cat a photo?*
- *For what reasons can I know that there is a cat in a photo?*

# Thiết kế thử nghiệm

**Experimental Design**

*(i.e. planning)*

**Question/ Hypothesis**

**Machine Learning**
*i.e.*
**Automatic data-driven predictive models**

Data? Acquisition?
*keywords: data sampling/survey*

Model? Assessment?
*keywords: training/testing sets, mean squared errors, precision, recall, …*

Interpre...

ssess...

...cquis...

Modelling

Pre-proces... (ETL)

**Thiết kế thử nghiệm**

Question/Hypothesis

Experimental Design
*(i.e. planning)*

Data sampling

Pre-process (ETL)

Modelling

Assessment

Interpret

**Machine Learning**
*i.e.*
**Automatic data-driven predictive models**

Data? Acquisition?
*keywords: data sampling/survey*

Model? Assessment?
*keywords: training/testing sets, evaluation metrics (e.g. mean squared errors, precision, recall)*

Question/
Hypothesis

Experimental
Design

## Data Sampling

### *Representative sample*

- *How many photos, categories, photos in each category, …?*
- *(If time-series data: eg videos)*
*Sample at which time points?*
- *Imbalance class?*
- *Selection bias?*

**Lấy mẫu**

Data sampling

Assessment

Modelling

Pre-pro
(ETL)

10

*Which metrics to use depend on which problem*
*http://scikit-learn.org/stable/modules/model_evaluation.html*

**Đánh giá mô hình**

## Model Assessment

**Evaluation metrics**

- *Accuracy*
- *Precision, Recall*
- *Area Under Curve (AUC)*
- *Mean squared errors (MSE)*
- *…*

*(If hypothesis testing problem)*

- *t-statistic, z-statistic, $\chi^2$-statistic, …*

Question/
Hypothesis

Experimental
Design

Data sampling

Assessment

Modelling

```
cat
flower
dog
jet
ground
grass
```

*If training/testing set split is well designed with sufficient examples, we might not need to repeat many experiments.*

Question/
Hypothesis

Experimental
Design

## Model Assessment

**Evaluation setup**

Evaluation (i.e.report results) on *unseen* data

- Training/testing set split: follows data sampling principles
- Repeat experiment: gives measurable confidence to the reported results

# Đánh giá mô hình

Assessment

Data sampling

Pre-pro
(ETL)

cat
flower
dog
jet
ground
grass

Modelling

*"All models are wrong, but some are useful."*
*- Box and Drape, 1987*

Question/
Hypothesis

## Model Building

Interpret

Experimental
Design

Model = a simplification of reality

*(e.g. map of Hanoi)*
*Keywords: Linear models, Graphical models, Neural networks,*
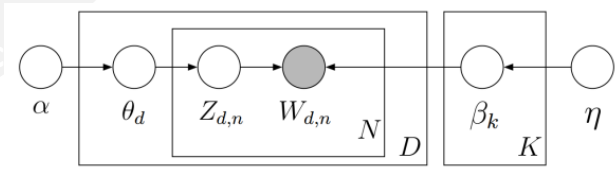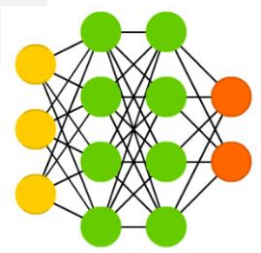*SVM, Gaussian Process, Random forest …*

Assessment

*Modelling tip*: building model goes from the <u>most</u>
<u>simplified</u> forms to the <u>more complex</u> to describe
reality more precisely
*(e.g. building from Linear models to Latent variable models /*
*Deep neural networks)*

Data acquisition

## Xây dựng mô hình

What ML
mostly
about

Modelling

Pre-proc
(ETL)

$\alpha$ → $\theta_d$ → $Z_{d,n}$ → $W_{d,n}$ $N$ $D$ | $\beta_k$ ← $\eta$ $K$

# Raw data ⟶ Post-processed data



- *Data ETL: extract, transform, load*
- *Data standardisation / normalisation*
- *Data imputation (if missing values)*

```
-0.34 -0.46 -0.87
 1.47 -0.24  2.21
-1.05  0.02 -1.74
 0.09 -0.58  1.02
 1.63 -0.53  0.06
 1.11 -0.63 -0.93
-0.34 -0.46 -0.87
 1.47 -0.24  2.21
-1.05  0.02 -1.74
 0.09 -0.58  1.02
 1.63 -0.53  0.06
 1.11 -0.63 -0.93
 0.09 -0.58  1.02
 1.63 -0.53  0.06
 1.11 -0.63 -0.93
 ....  ....  ....
```
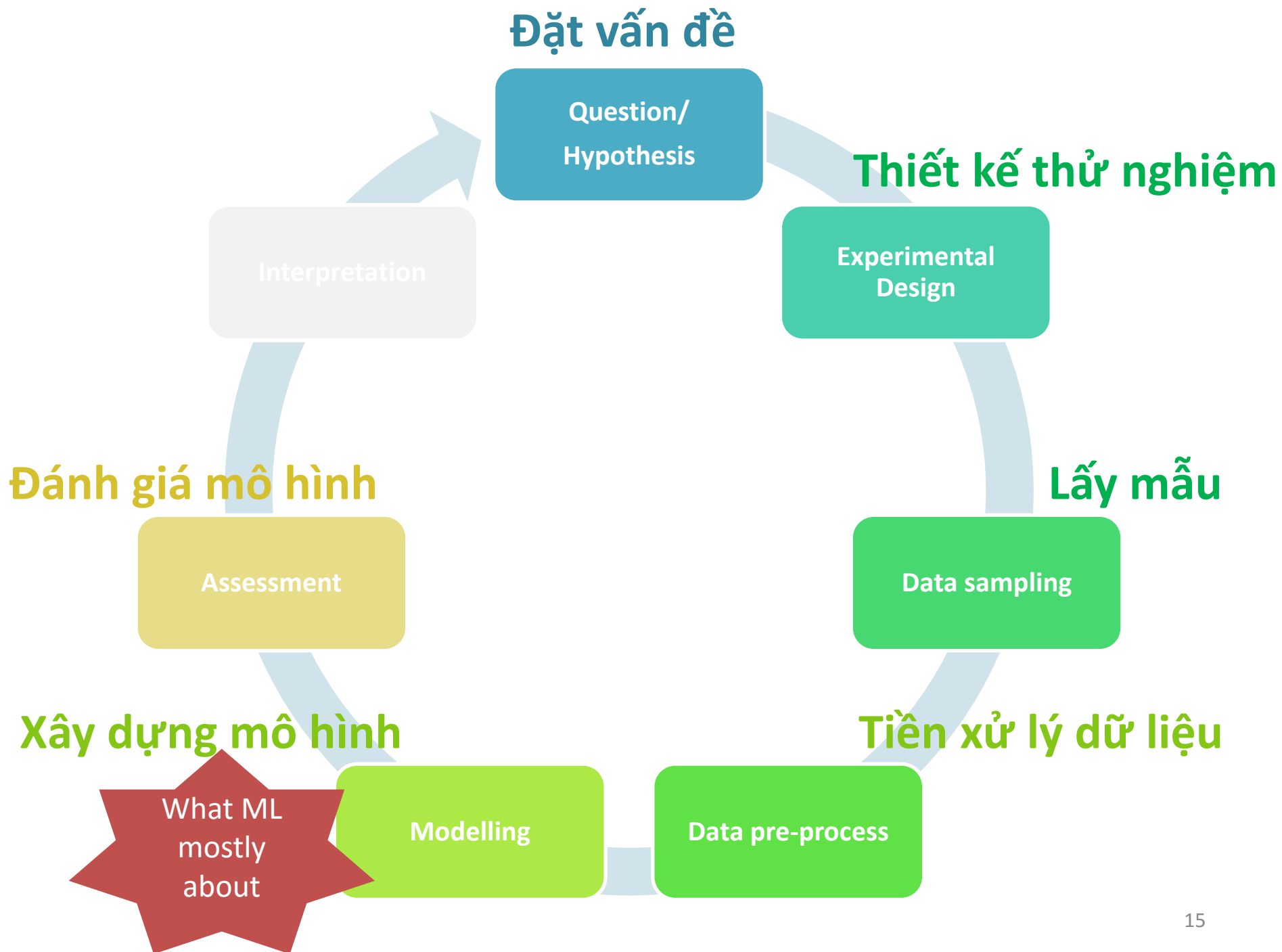
*Feature extraction*

*Foreshadowing*: the core idea of **Deep Learning** is to incorporate *feature extraction* stage into a model, for which how the features are extracted is also *learnt from the data*.

**Tiền xử lý dữ liệu**

**Modelling**  **Data pre-process**

Vấn đề, câu hỏi mới

NEW Question/ Hypothesis

Thiết kế thử nghiệm

Experimental Design

Giải thích/phân tích kết quả

Interpretation

Lấy mẫu

Data sampling

Đánh giá mô hình

Assessment

Tiền xử lý dữ liệu

Data pre-process

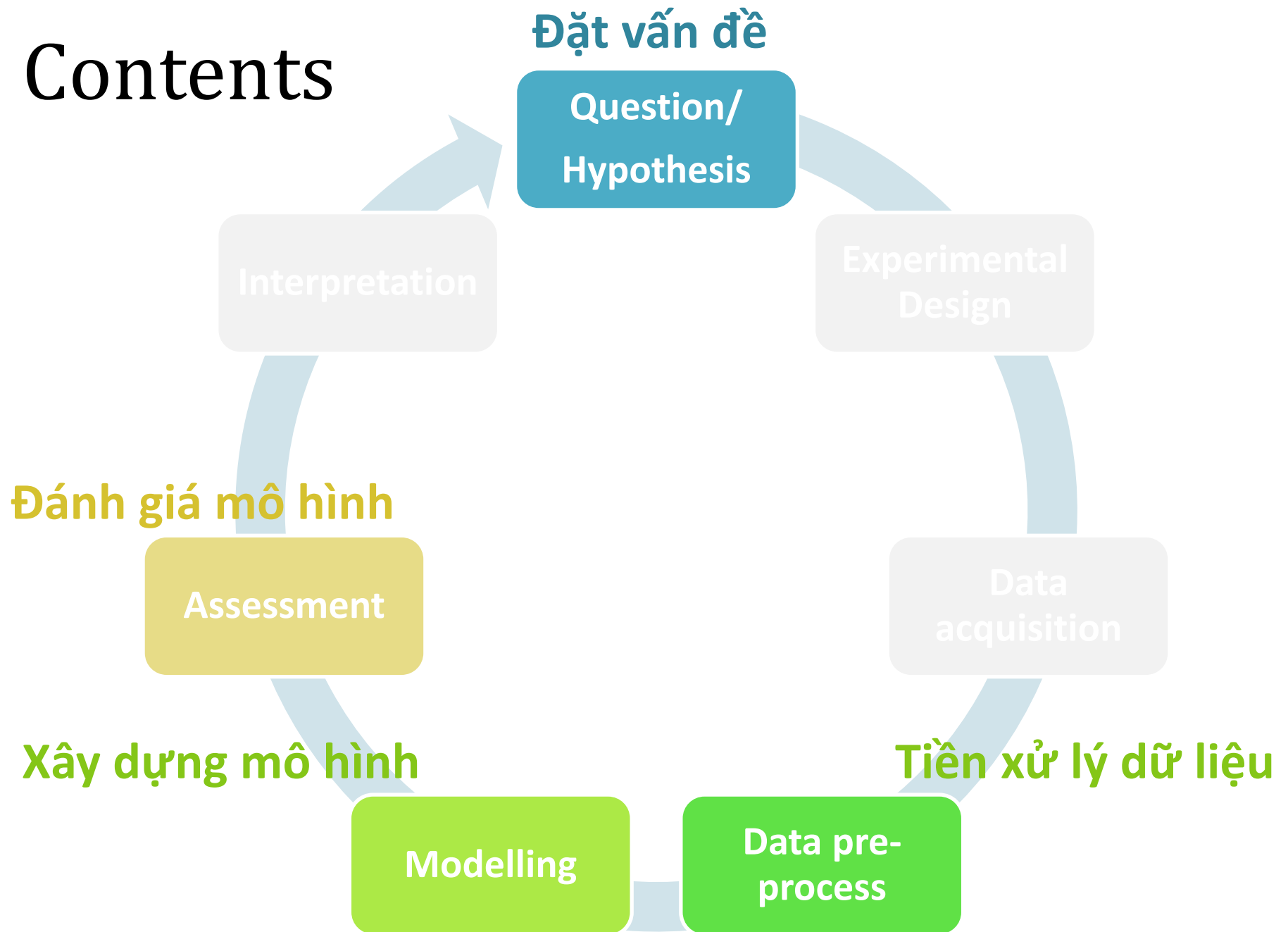Xây dựng mô hình

Modelling

What ML mostly about

# PRINCIPLES OF MODELLING

## Statistical reasoning $^{(*)}$

*(*) A machine learning algorithm does not necessarily have a probabilistic interpretation, or developed from a statistical framework. Nevertheless, statistical reasoning provides a rigorous mathematical tool for estimation and inference to make optimal decision (e.g. prediction, action) under **uncertainty,** which is one of the ultimate objectives in ML.*

# Contents

# ML problem: Classification

**Question**  **Is there** any cat in an abitrary photo?

Experience: dataset of {image, label} pairs $\mathcal{D} = \{x_n, y_n\}_{n=1}^{N}$

**Modelling**  predict $\hat{y}_n$ – *cat existence* – *given arbitrary* $x_n$



Cat?
Not cat?

**supervised learning**

**Image**
$x_n$
$\mathbb{N}^{400 \times 600 \times 3}$

**Prediction**
$\hat{y}_n$
{True, False}

*(single-class)*

**binary classification problem**

**Assessment**  Accuracy $= \frac{1}{N} \sum_n \mathbb{I}(\hat{y}_n = y_n)$

Precision, Recall, F1-score

Area Under Curve (AUC)

...

*Example models:*
*Logistic regression (linear model)*
*Neural Net with sigmoid output (nonlinear model)*

# ML problem: Classification

**Question**

**What is there in an abitrary photo?**

Experience: dataset of {image, label} pairs $\mathcal{D} = \{x_n, y_n\}_{n=1}^{N}$

**Modelling**

predict $\hat{y}_n$ – *object identity* – *given arbitrary* $x_n$



**cat**
flower
dog
jet
ground
grass

**supervised learning**

**Image**
$x_n$
$\mathbb{N}^{400 \times 600 \times 3}$

**Prediction**
$\hat{y}_n$
$\{1,2,3,4,5,6\}$

*(multi-class)*

**categorical classification problem**

**Assessment**

Accuracy $= \frac{1}{N} \sum_n \mathbb{I}(\hat{y}_n = y_n)$

Precision, Recall, F1-score

Area Under Curve (AUC)

...

*Example models:*
*Softmax classification (linear model)*
*Neural Net with softmax output (nonlinear model)*

# ML problem: Regression

**Question** **How much** is the price of a house given …

Experience: dataset of {(area, location, #rooms), price} pairs $\mathcal{D} = \{x_n, y_n\}_{n=1}^{N}$

**Modelling** predict $\hat{y}_n$ – *house price* – *given arbitrary* $x_n$

| Area | 100m$^2$ |
|------|----------|
| Location | 24.7$^0$N 183.0$^0$E |
| #Rooms | 3 |

→ $150,000

**supervised learning**

**Features/Predictors**
$$x_n$$
$$\mathbb{R} \times \mathbb{R}^2 \times \mathbb{N}$$

**Prediction**
$$\hat{y}_n$$
$$\mathbb{R}$$

**regression problem**

**Assessment** squared_errors = $\frac{1}{N}\sum_n(\hat{y}_n - y_n)^2$

*Example models/algorithms:*
*Linear regression (linear model)*
*Neural Net with linear output (nonlinear model)*
*Curve fitting algorithm*

# ML problem: Clustering

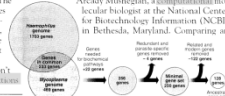**Question** — **What is** the "topic" that a news article is talking about?

Experience: dataset of article content $only \; \mathcal{D} = \{x_n\}_{n=1}^{N}$

**Modelling** — predict $\boxed{z_n}$ – *"topic" (cluster) identity* – *given arbitrary* $x_n$



**unsupervised learning**

Graphics:
- David Blei, KDD 2011
- https://lvdmaaten.github.io/tsne/examples/mnist_tsne.jpg

**Article (text)**
$x_n$
$\mathbb{N}^{1500}$

**Prediction**
$z_n$
$\{1, 2, \dots, 10\}$

$x_n$
$z_n = $ **green**

**Assessment** — $\text{mean\_distance\_to\_clusters} = \frac{1}{N} \sum_n \left( x_n - \mu_{z_n} \right)^2$

*Note: "topic" = group/cluster in this context, and is <u>not</u> pre-defined*
*We will meet the term "topic" again when visiting Topic models*

*Example models/algorithms:*
*k-means algorithm*
*Generative models: Mixture models, Topic models*

A ML problem can also be:

- both supervised and unsupervised *(semi-supervised)*
- combination of regression and classification sub-problems *e.g. image localisation*



**Classification**: C classes
  **Input**: Image
  **Output**: Class label
  **Evaluation metric:** Accuracy

→ CAT

**Localization**:
  **Input**: Image
  **Output**: Box in the image (x, y, w, h)
  **Evaluation metric:** Intersection over Union

→ (x, y, w, h)

**Classification + Localization**

CAT

# PRINCIPLES OF MODELLING

1.  **Model structure -** constructs relationships (*stochastic and/or deterministic*) between model elements: data, parameters, and hyper-parameters.

    *Keywords: graphical model*

2.  **Learning principle -** defines a framework to estimate unknown parameters (and unobserved i.e. hidden/latent variables)

    *Keywords: Maximum Likelihood criterion, Bayesian inference, ++ others*

3.  **Regularisation**

    *Keywords: over-fitting, Bayesian inference, ++ others*
    *Relevant keywords: L2-regularisation (Ridge), L1-regularisation (LASSO)*

⇒ **ALGORITHM -** implements 1 + 2 + 3 to train the model

    *Keywords: (stochastic) gradient descent, Expectation-Maximisation (EM), Variational Inference (VI), sampling-based inference methods*

4.  **Model selection**

    *Keywords: cross-validation*

Before we get going…

"Mathematics is the art of giving the same name to different things ."

-Henri Poincaré.

"The purpose of computation is insight, not numbers."

-Richard Hamming

$$p(\mathbf{w}\,|\,\alpha,\beta) = \frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \int \left( \prod_{i=1}^{k} \theta_i^{\alpha_i - 1} \right) \left( \prod_{n=1}^{N} \sum_{i=1}^{k} \prod_{j=1}^{V} (\theta_i \beta_{ij})^{w_n^j} \right) d\theta,$$

$$p(D\,|\,\alpha,\beta) = \prod_{d=1}^{M} \int p(\theta_d\,|\,\alpha) \left( \prod_{n=1}^{N_d} \sum_{z_{dn}} p(z_{dn}\,|\,\theta_d) p(w_{dn}\,|\,z_{dn},\beta) \right) d\theta_d.$$