# 2016 NECINA Financial Technology Conference

## Applying Machine Intelligence in Financial Software Engineering
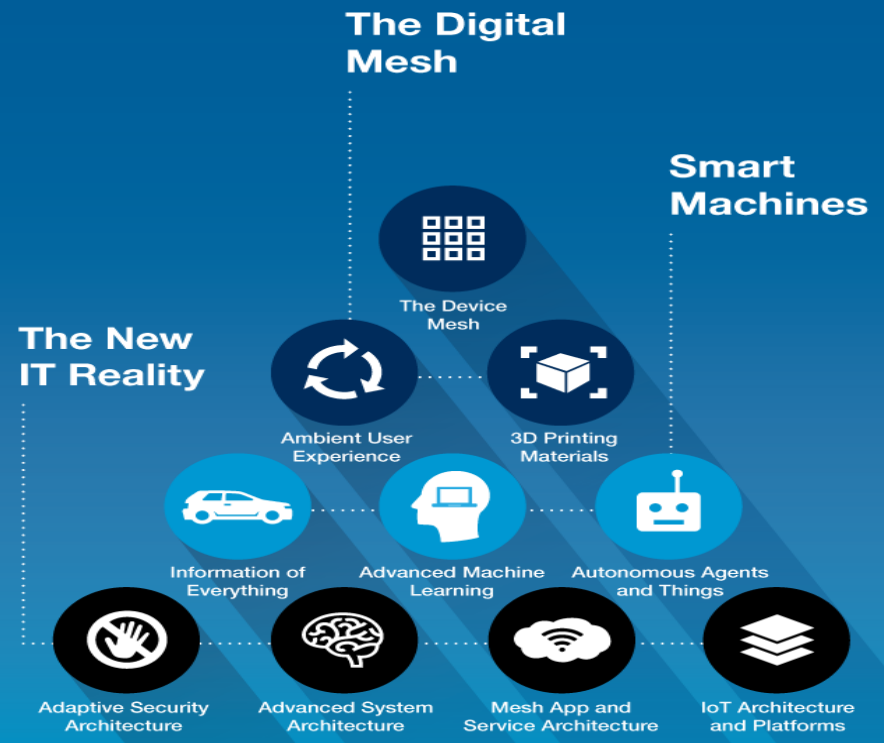
By Albert Ma, Chief Innovation Officer of Hengtian

Date: April 23th, 2016

**HengTian**
Trusted Technology Solutions

# Agenda

- Machine Intelligence Landscape
- Business Values
- Auto Programming
- DeepMorpho – Research for the Future
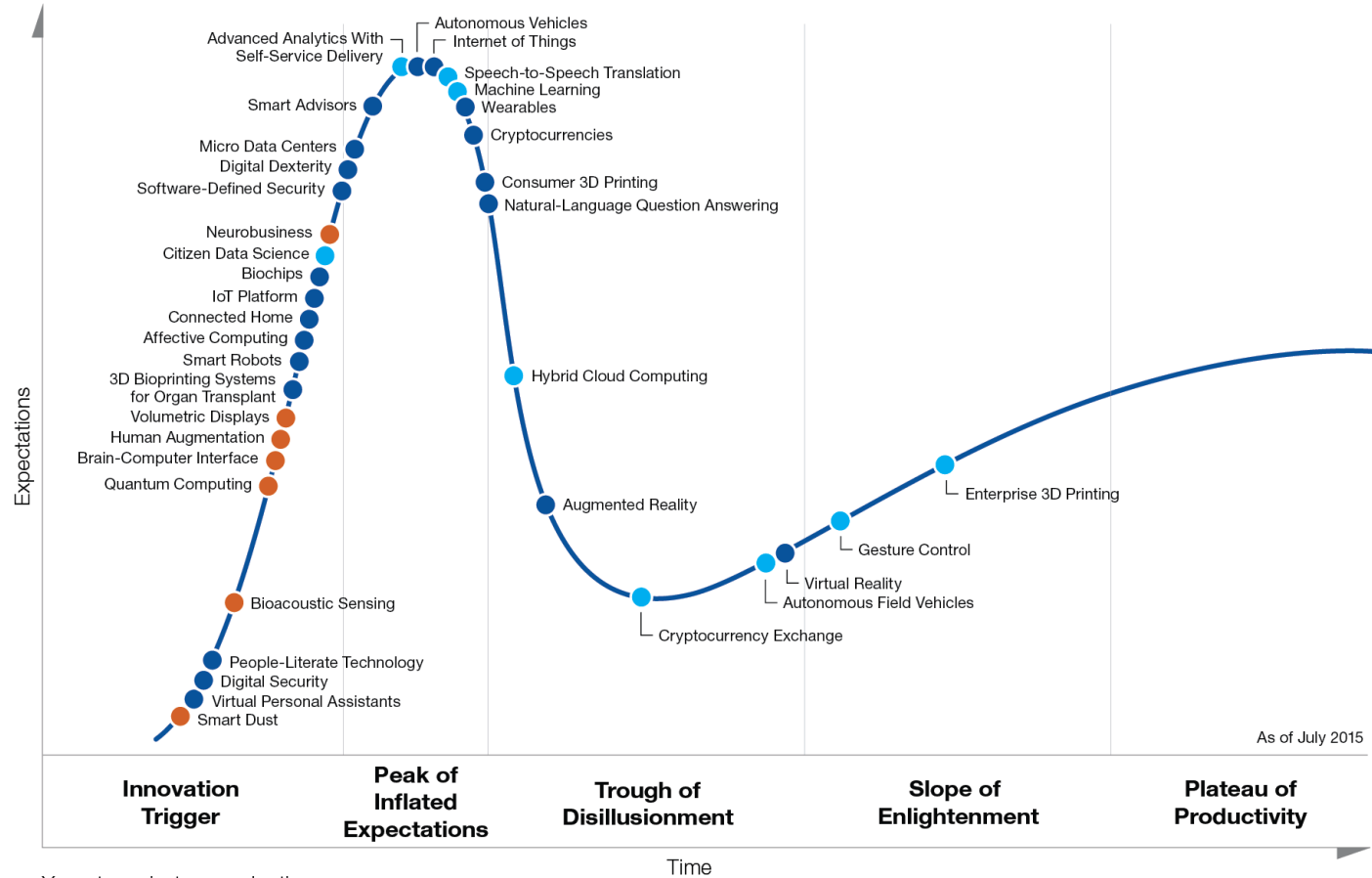- Summary

HengTian
Trusted Technology Solutions

# Emerging Technology Hype Cycle



**Expectations**

Advanced Analytics With Self-Service Delivery
Autonomous Vehicles
Internet of Things
Speech-to-Speech Translation
Machine Learning
Wearables
Smart Advisors
Cryptocurrencies
Micro Data Centers
Digital Dexterity
Software-Defined Security
Consumer 3D Printing
Natural-Language Question Answering
Neurobusiness
Citizen Data Science
Biochips
IoT Platform
Connected Home
Affective Computing
Smart Robots
3D Bioprinting Systems for Organ Transplant
Volumetric Displays
Human Augmentation
Brain-Computer Interface
Quantum Computing
Hybrid Cloud Computing

Augmented Reality

Enterprise 3D Printing

Gesture Control

Virtual Reality
Autonomous Field Vehicles

Cryptocurrency Exchange

Bioacoustic Sensing

People-Literate Technology
Digital Security
Virtual Personal Assistants
Smart Dust

As of July 2015

| Innovation Trigger | Peak of Inflated Expectations | Trough of Disillusionment | Slope of Enlightenment | Plateau of Productivity |

**Time**

Years to mainstream adoption:

● less than 2 years  ● 2 to 5 years  ● 5 to 10 years  ● more than 10 years  ⊗ obsolete before plateau

**Gartner.**

**HengTian**
Trusted Technology Solutions

# The Trends of Machine Intelligence

▶ Two Drivers:  Raw Data and Data Model

▶ Machine Learning becomes a buzz word for business (InfoWorld)

▶ Big Data is assumed in Machine Learning applications (Gartner)

▶ Wall Street is gearing up with Machine Learning in fixed income, blockchain, predictive analytics etc. (McKinsey)

▶ Machine Learning is full of contradiction (Thomas Frey, The Da Vinci Inst.)

HengTian
Trusted Technology Solutions

# It is the year of artificial intelligence !

When DeepMind AlphaGo
win Lee Sedol…….

Can computer read source
codes like a human ?

?

```
COBOL Code
DETERMINE-PMHP.                                                    04380074
043900   IF PMHP-FAC NOT = CC-FAC                                  04390074
044000      MOVE 'Y' TO EOF-PMHP                                   04400074
044100      GO TO EXIT-PMHP.                                       04410074
044200   INITIALIZE SORT-RECORD.                                   04420074
044300   MOVE PMHP-FAC  TO SORT-FAC.                               04430074
044400   MOVE PMHP-CASE TO SORT-CASE.                              04440074
044500   IF PMHP-INELIGIBLE-CODE NOT = ZERO                        04450074
044600      MOVE 5 TO SORT-PMHP-STATUS                             04460074
044700   ELSE IF PMHP-ACCEPT-DECLINE-FLAG = 2                      04470074
044800      MOVE 4 TO SORT-PMHP-STATUS                             04480074
044900   ELSE IF (PMHP-ENROLL-DSS-RESPONSE = 01 OR 02)             04490074
045000      AND PMHP-DISENROLL-DSS-DATE = ZEROES                   04500074
045100      MOVE 1 TO SORT-PMHP-STATUS                             04510074
045200   ELSE IF PMHP-DISENROLL-DSS-DATE NOT = ZERO                04520074
045300      MOVE 3 TO SORT-PMHP-STATUS                             04530074
045400   ELSE IF PMHP-ENROLL-EXTRACT-DATE = ZEROES                 04540074
045500      AND PMHP-DISENROLL-REASON NOT = ZERO                   04550074
045600      MOVE 3 TO SORT-PMHP-STATUS                             04560074
045700   ELSE                                                      04570074
045800      MOVE 2 TO SORT-PMHP-STATUS                             04580074
045900      IF PMHP-CORRECTION-DATE > PMHP-ENROLL-EXTRACT-DATE     04590074
046000         MOVE 1 TO SORT-READY-RESEND
```
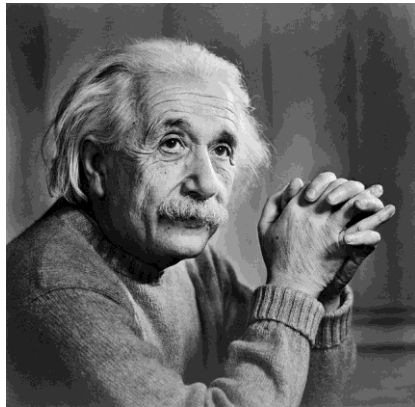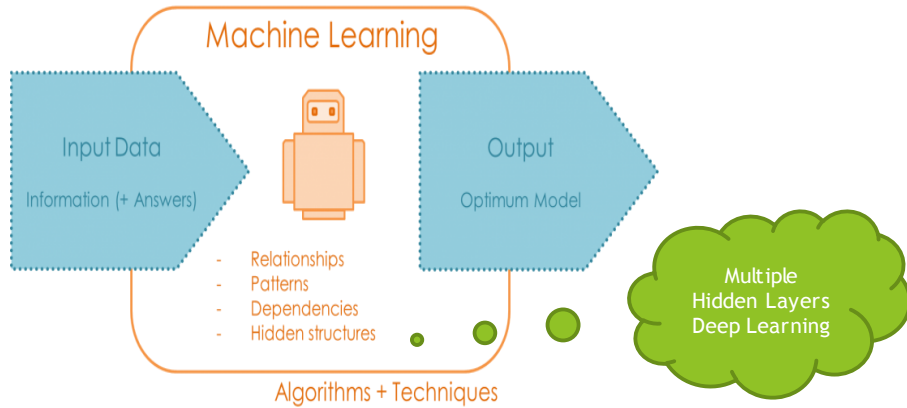
**HengTian**
Trusted Technology Solutions

# Auto Programming Samples

- Tweet to program (http://www.wolfram.com)

- Natural Language Programming (http://www.pegasus-project.org)

- Cognitive Computing (http://www.ibm.com/watson)

- Statistical Machine Translation (https://youtube.com/watch?v=aRSnl5-7vNo)

- Pliny Big Code Analytics (http://pliny.rice.edu/index.html)

- Code Transplant ( http://crest.cs.ucl.ac.uk/autotransplantation/MuScalpel.html)

- Bug Repair - MIT CodePhage (http://news.mit.edu/2015/automatic-code-bug-repair-0629)

- Malware Prevention - Deep Instinct  ( http://www.deepinstinct.com/#/what-we-do )

HengTian
Trusted Technology Solutions

# Some Basics of Machine Learning

## Machine Learning

Input Data
Information (+ Answers)

Output
Optimum Model

- Relationships
- Patterns
- Dependencies
- Hidden structures

Algorithms + Techniques

Multiple Hidden Layers Deep Learning

### Supervised Learning:

Predicting values. **Known** targets.
User inputs correct answers to learn from. Machine uses the information to guess new answers.

**REGRESSION**:
Estimate continuous values
(Real-valued output)

**CLASSIFICATION**:
Identify a unique class
(Discrete values, Boolean, Categories)

### Unsupervised Learning:

Search for structure in data. **Unknown** targets.
User inputs data with undefined answers. Machine finds useful information hidden in data.

**Cluster Analysis**
Group into sets

**Density Estimation**
Approximate distributions

**Dimension Reduction**
Select relevant variables

## Supervised Learning

### Regression
- Linear Regression
- Ordinary Least Squares Regression
- LOESS (Local Regression)
- Neural Networks

### Classification
- Decision Trees
- K-Nearest Neighbours
- Support Vector Machine
- Logistic Regression
- Naïve Bayes
- Random Forests

## Unsupervised Learning

### Cluster Analysis
- K-Means Clustering
- Hierarchical Clustering

### Dimension Reduction
- Principal Component Analysis (PCA)
- Linear Discriminant Analysis (LDA)

### *TRAINING*
Learn data properties

The machine makes conclusions by learning from the data.

It improves its model until optimal performance is reached.

Using a Cost / Loss Function to measure accuracy. It repeats iterations until a minimum is reached. (e.g. gradient descent)

### *TESTING*
Test the properties

Apply the conclusions to new data and compare results to known answers.

The model does not change. It is just tested to measure how good the machine did after training.

Useful to detect overfitting. If good enough, it is ready to be used.

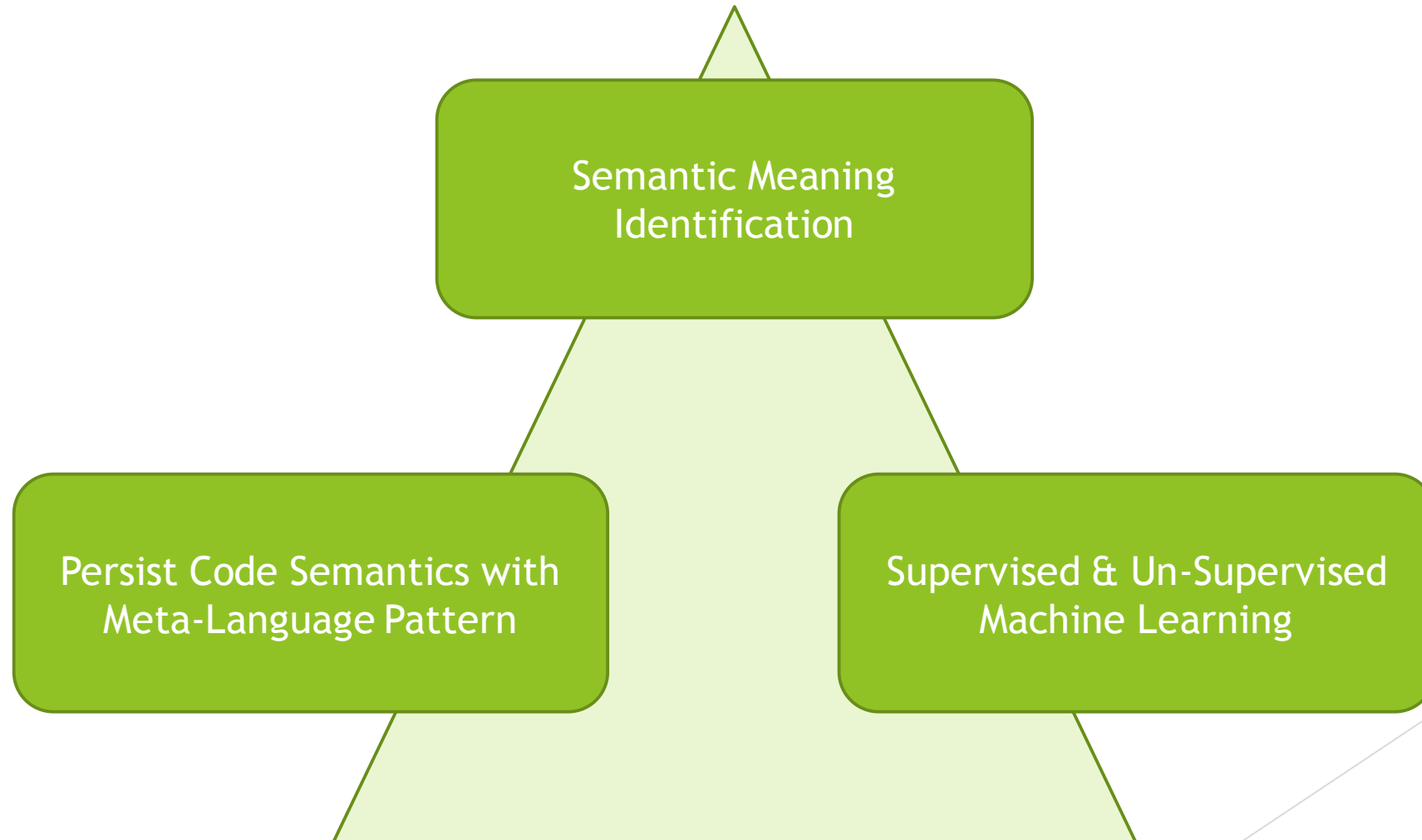### *APPLICATION*
Use the properties

In a real situation the answers are not known.

Apply the model conclusions to predict the answers from the inputs. Use the answers in whatever necessary.

Source: http://quantdare.com/2016/03/machine-learning-a-brief-breakdown/

**HengTian**
Trusted Technology Solutions
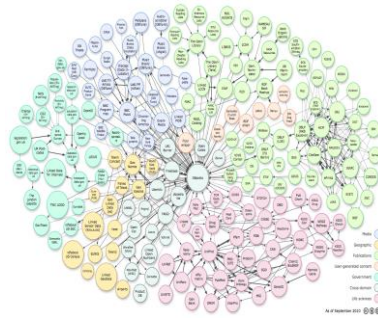
# Why learning from source code is difficult ?

▶ Programing language syntax structure is well defined in the complier but has infinite ways of implementing the same business logic

▶ There is not a universal way of representing program semantic meanings

▶ Existing NLP representation learning algorithms are inapplicable since all of them are "flat"

▶ Program symbols (nodes in AST) are discrete and cannot be fed directly to a neural network

▶ Like NLP, it takes many years to mature multi-language corpus

▶ Multiple programing languages and systems interconnected together inflates the complexity permutation

**HengTian**
Trusted Technology Solutions

# The problems to be resolved in auto programming?

Semantic Meaning Identification

Persist Code Semantics with Meta-Language Pattern

Supervised & Un-Supervised Machine Learning

HengTian
Trusted Technology Solutions

# Value Proposition – Why Machine Learning ?

**1**



Enterprise Code Ontology

Know your IT asset at granular level

**2**

Malware ?

Bugs ?
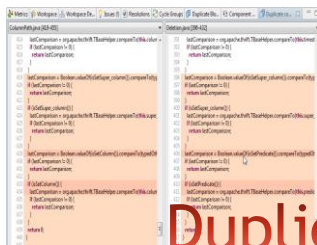
System Failure ?

Back Door ?



Predictive Analytics

**3**



Duplicate & Dead Codes

**4**



Code Transplant

HengTian
Trusted Technology Solutions
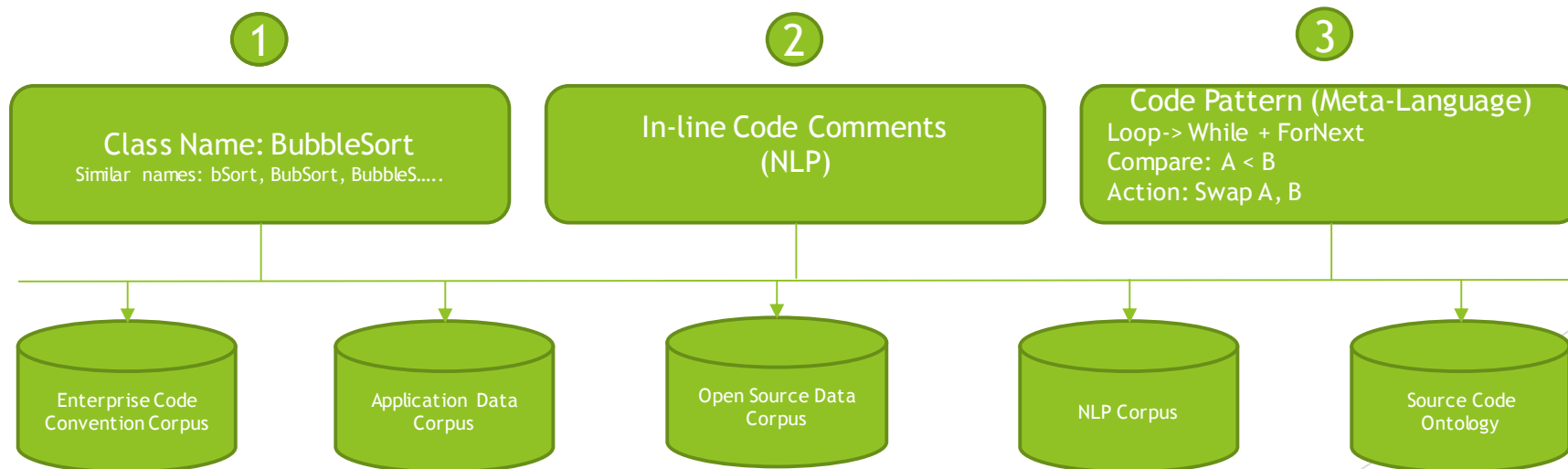
# How does a programmer read the code ?

```
public static void BubbleSort( int [ ] num ) {
int j;
boolean sortFlag = true;  // initialize the flag for while loop
    int tempNum;

while ( sortFlag ) {

        sortFlag = false;    //set flag to false awaiting a possible swap
        for( j=0;  j < num.length -1;  j++ ){
                if ( num[ j ] < num[j+1] ) {  // compare adjacent array elements

                        tempNum = num[ j ];          //swap two elements
                        num[ j ] = num[ j+1 ];
                        num[ j+1 ] = tempNum;
                        sortFlag = true;            //shows a swap occurred
                }
        }
    }
}
```

Bubble Sort ?

| ① | ② | ③ |
|---|---|---|
| **Class Name: BubbleSort**<br>Similar names: bSort, BubSort, BubbleS..... | **In-line Code Comments<br>(NLP)** | **Code Pattern (Meta-Language)**<br>Loop-> While + ForNext<br>Compare:  A < B<br>Action: Swap A, B |

Enterprise Code Convention Corpus

Application Data Corpus

Open Source Data Corpus

NLP Corpus

Source Code Ontology

**Code Signature**

HengTian
Trusted Technology Solutions

# Source Code Representation

- Semantic Annotation
  - Triples ('class A', 'inherit from', 'class B'); ('class A', 'has bizfunc', 'interest calculation')
  - Resource Description File (RDF)
- Ontological Representation for Machine Learning
  - Vector based code block
  - Data flow
- Meta Language for Pattern Persistence
  - Regular Expression
  - Domain Specific Language (DSL)

# Summary

▶ There is not a single model that can solve all the software development problems

▶ Significant effort is needed to develop different language corpus and training data set. Payback occurs when there is a strategic motive and sustainable methods embedded in software development life cycle

▶ It is just the beginning for academic researchers and software vendors to look at applying machine intelligence into software development. The potential is huge with more and more algorithmic libraries coming out

▶ Accuracy will significantly improve if there is an enterprise code standard and it is being enforced properly

▶ As more software development counts on open source, it is a logical next step to extend machine intelligence to transplant codes from open source repositories

**HengTian**
Trusted Technology Solutions

# Recommended Papers

- L.Mou, G.Li, Y.Liu, H.Peng, Z.Jin, Y.Xu, L.Zhang, Building Program Vector Representations for Deep Learning, Software Institute, School of EECS, Peking University

- N.Phan, D.Dou, H.Wang, D.Kil, Ontology-Based Deep Learning for Human Behavior Prediction in Health Social Networks, Computer Science, University of Oregon

- F.Long, M.Rinard, Automatic Patch Generation by Learning Correct Code, CSAIL, Massachusetts Institute of Technology

- E.Barr, M.Harman, Y.Jia, A.Marginean, J.Petke, Automated Software Transplantation, CREST, University College London

- W.Ling, E.Grefenstette, K.M.Hermann, T.Kocisky, A.Senior, F.Wang, P.Blunsom, Latent Predictor Networks for Code Generation, Google DeepMind, University of Oxford

- S. Liu, N.Yang, M.Li, M.Zhou, A Recursive Recurrent Neural Network for Statistical Machine Translation, University of Science and Technology of China