# UNSUPERVISED HANDWRITTEN GRAPHICAL SYMBOL LEARNING
# -Using Minimum Description Length Principle on Relational Graph

Jinpeng LI, Harold MOUCHERE, Christian VIARD-GAUDIN
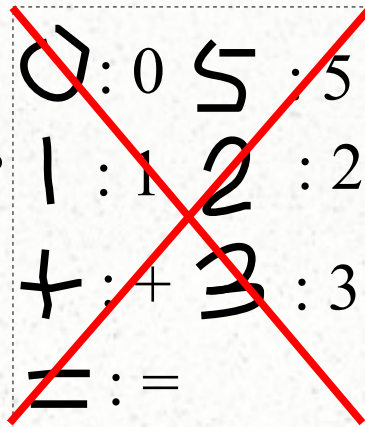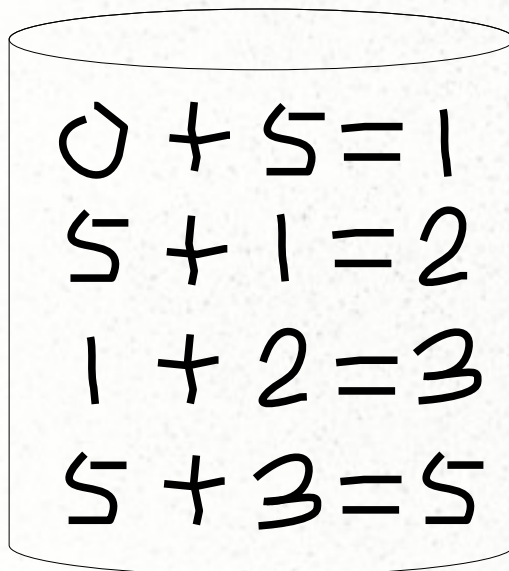
KDIR 2011
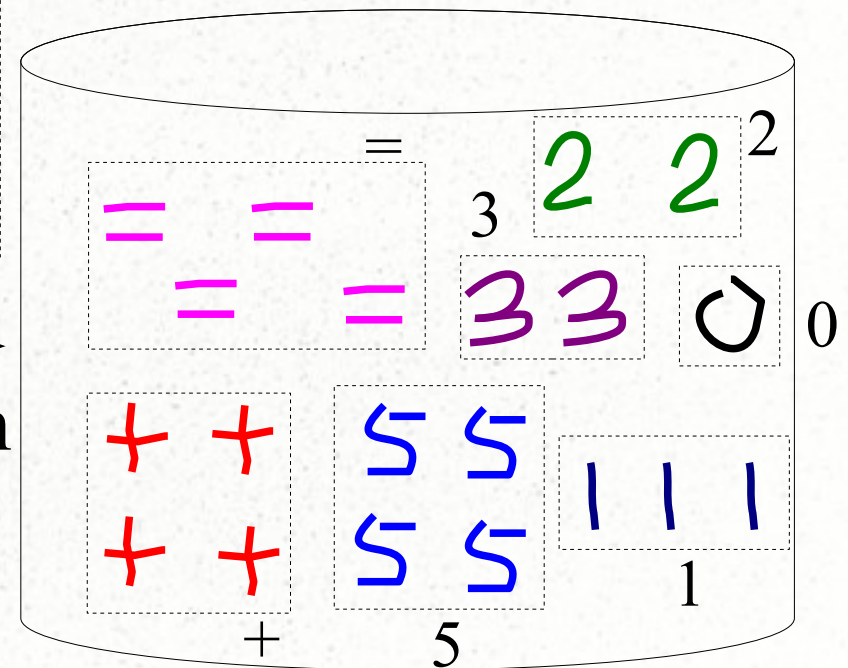
www.projet-depart.org

# *Outline*

- 1. Background
- 2. Unsupervised Handwritten Graphical Symbol Learning
    - Quantization of strokes
    - Relational Graph Construction Between Strokes
    - Discover Symbols (Sub-graphs)
- 3. Experiment
- 4. Conclusion

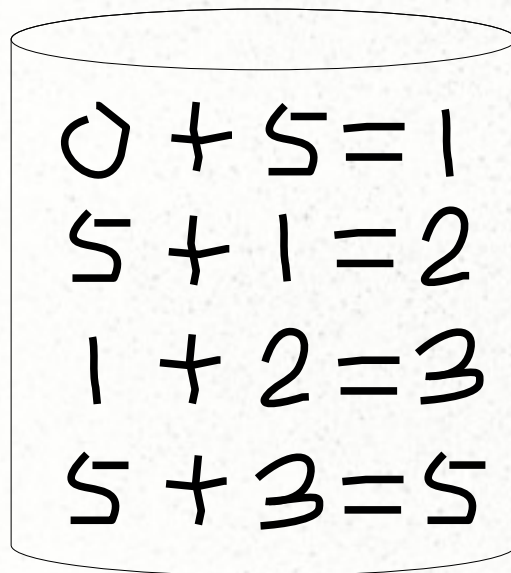# *Traditional Recognition (Background)*
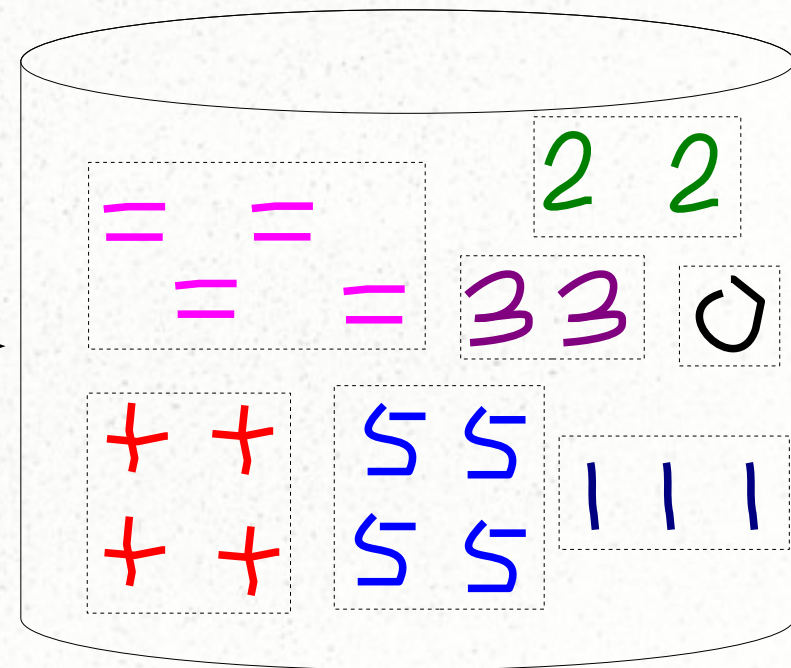
Unlabeled handwritten symbols



Recognition

# *Annotation (Application)*

Unlabeled
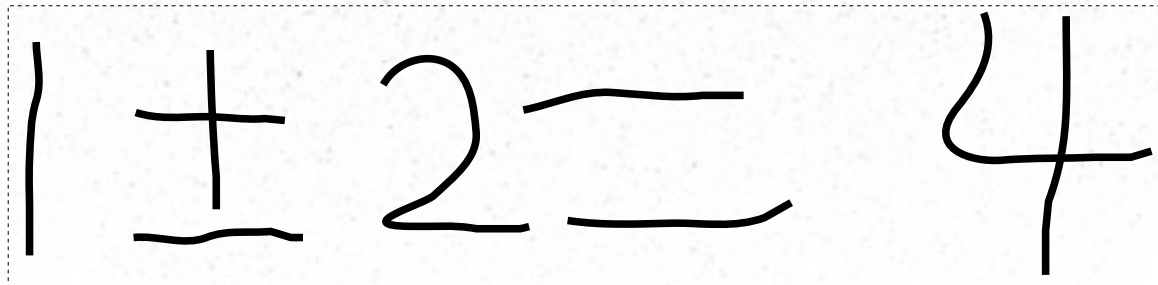handwritten symbols



Symbol
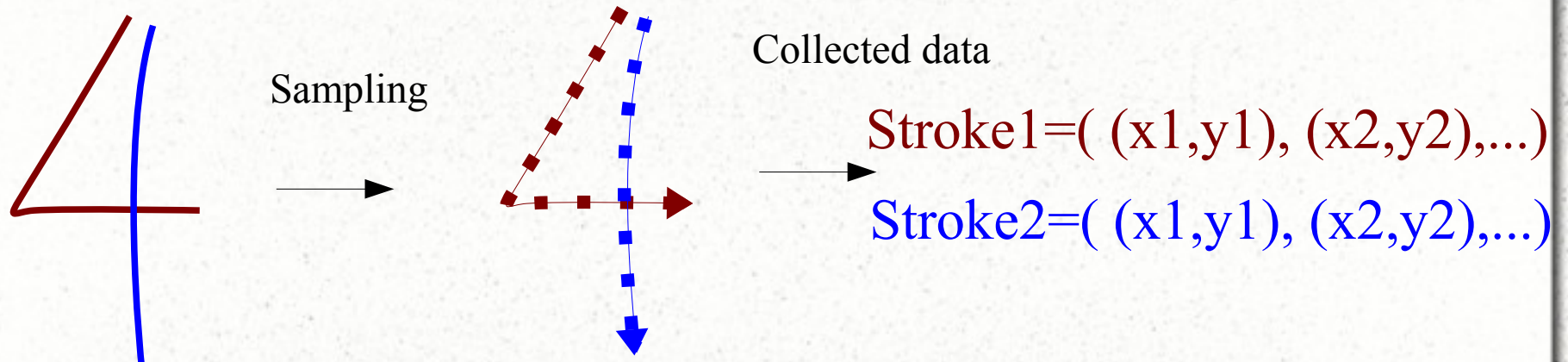Extraction

**40** symbols
have to be labeled

**7** symbols (sets)
have to be labeled

# UNSUPERVISED HANDWRITTEN GRAPHICAL SYMBOL LEARNING

As an example, we use mathematical expressions as an **unknown** graphical language.

# *Online handwriting*
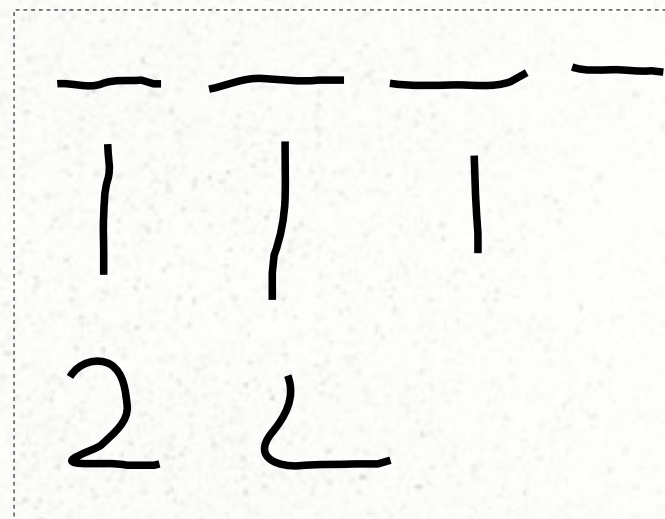
Sampling

Collected data

Stroke1=( (x1,y1), (x2,y2),...)

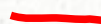Stroke2=( (x1,y1), (x2,y2),...)

# *Unsupervised symbol learning*

$1 \pm 2 = 4$ $\xrightarrow{\text{The base elements are strokes}}$

We call the frequent strokes as
the **Grapheme**.

# *Unsupervised symbol learning*

$1 \pm 2 = 4$  The base elements are strokes →

The horizontal stroke repeats 4 times.

We call the frequent strokes as the **Grapheme**.

# *Grapheme*

$1 \pm 2 = 4$

From a part of symbol "plus".

$1 \pm 2 = 4$

From a symbol "equal".

$1 \pm 2 = 4$

From a symbol "minus".

Where does
the horizontal stroke
come from?

# *Outline*

- 1. Background
- 2. Unsupervised Handwritten Graphical Symbol Learning
    - **Quantization of strokes**
    - Relational Graph Construction Between Strokes
    - Discover Symbols (Sub-graphs)
- 3. Experiment
- 4. Conclusion

# *Hierarchical clustering*



Strokes

Dynamic Time Warping Distance

# *Hierarchical clustering*



Strokes

Codebook

Dynamic Time Warping Distance

12

# *Quantization of strokes*
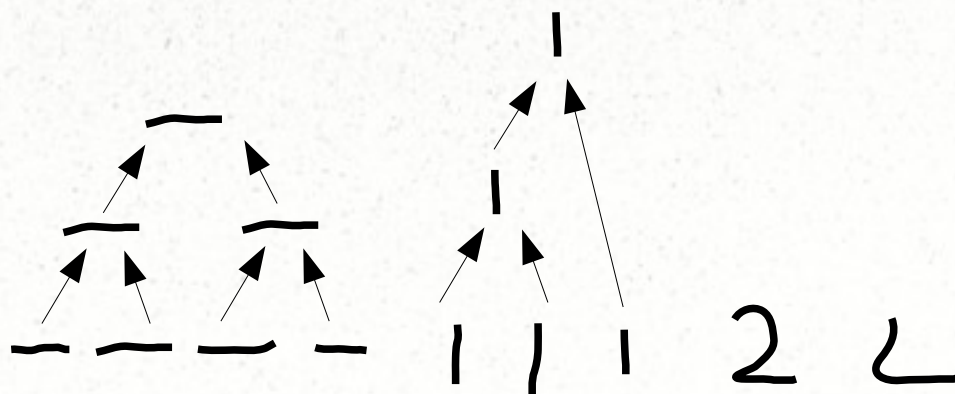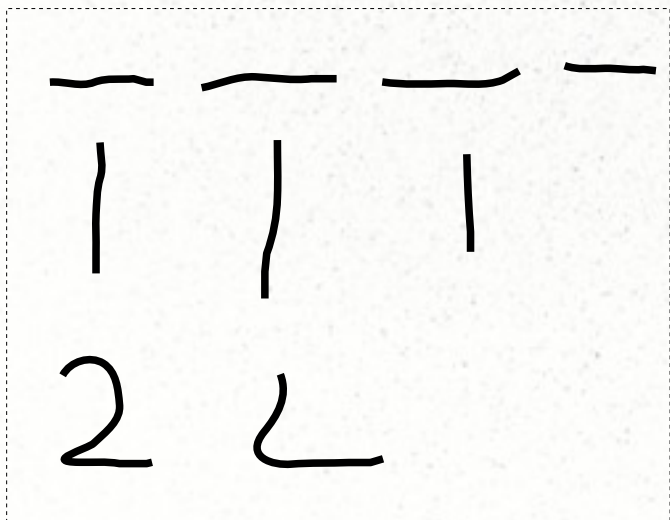
Codebook

1. Quantization of strokes

# *Outline*

- 1. Background
- 2. Unsupervised Handwritten Graphical Symbol Learning
    - Quantization of strokes
    - **Relational Graph Construction Between Strokes**
    - Discover Symbols (Sub-graphs)
- 3. Experiment
- 4. Conclusion

# *Relational Graph Construction*

Spatial relation: from a reference stroke to an argument stroke.

Node

Stroke Ref

Edge: Spatial relation

Stroke Arg

# *Relational Graph Construction*

Spatial relation: from a reference stroke to an argument stroke.

We predefine three spatial relations:
right ( R ), below ( B ), and intersection ( I ).



$n_{str}$ : the number of strokes (number of nodes).

$n_r$ : the number of different spatial relations from a reference stroke to an argument stroke.

$n_r = 3$

# *Relational Graph Construction*

Spatial relation: from a reference stroke to an argument stroke.

We predefine three spatial relations:
right ( R ), below ( B ), and intersection ( I ).

We predefine another constraint that
Directional spatial relation (R and B )
are exclusive with
Topological spatial relation (I).

Stroke Ref          Stroke Ref

Right          Below   **Or**          Intersection

Stroke Arg          Stroke Arg

$n_{str}$ : the number of strokes (number of nodes).

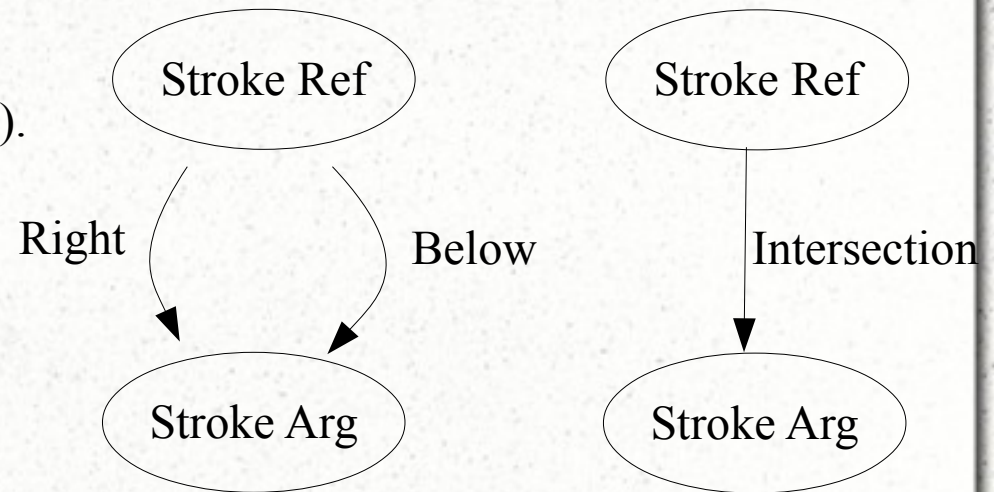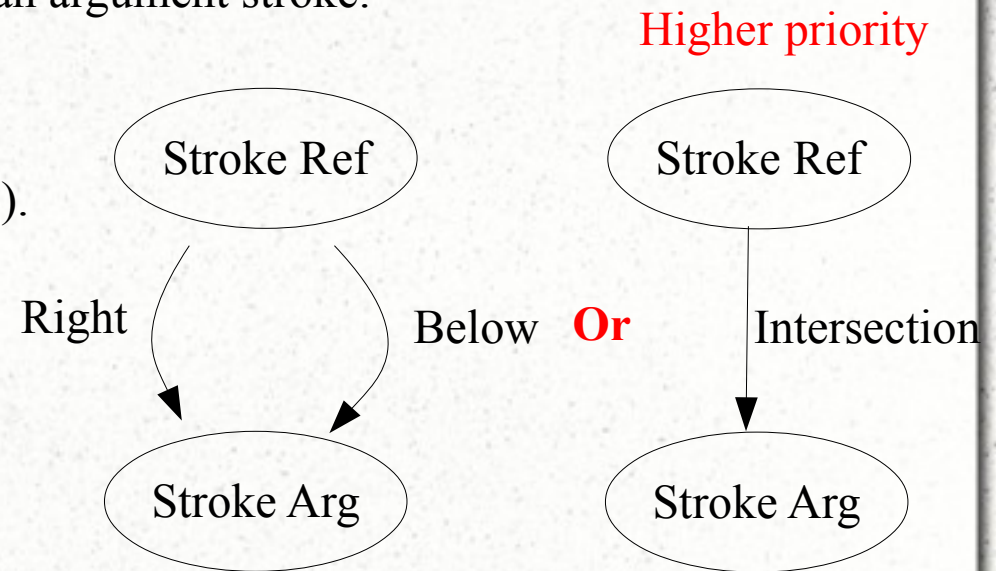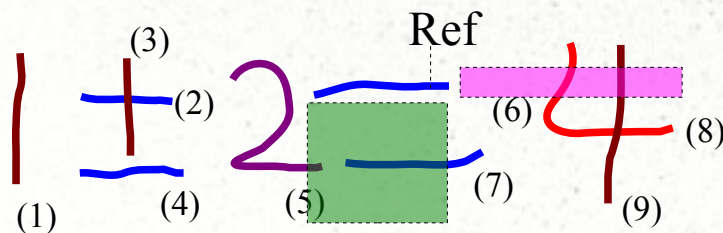$n_r$ : the number of different spatial relations from a reference stroke to an argument stroke.

$n_r = 2$

17

# *Relational Graph Construction*

Spatial relation: from a reference stroke to an argument stroke.

Higher priority

We predefine three spatial relations:
right ( R ), below ( B ), and intersection ( I ).

We predefine another constraint that
Directional spatial relation (R and B )
are exclusive with
Topological spatial relation (I).

Stroke Ref

Right     Below   **Or**

Stroke Arg

Stroke Ref

Intersection

Stroke Arg

(3)

Ref

(2)

(6)

(8)

(1)     (4)     (5)     (7)     (9)

Right:

Below:

R

(6)          (8)

B     B     R

(5)     (7)     (9)     18

# *Number of edges*

Spatial relation: from a reference stroke to an argument stroke.

$n_r$ : the number of different spatial relations.

$n_{str}$ : the number of strokes.

(3) (2) (6) (8)
(1) (4) (5) (7) (9)

$n_r = 2$     $n_{str} = 9$

Complete directed graph

**Too many edges!**

Number of edges in graph

$O(n_{str}^2)$   $n_r n_{str}(n_{str} - 1)$ ........................................ $=144$

$n_c \leq (n_{str} - 1)$

We prefer some symbols
composed of the $n_c = 2$ closest strokes

**Reduced**

$n_r \cdot n_{str} \cdot n_c$ ........................................ $2 \cdot 9 \cdot 2 = 36$

Since we, human, have a limited perceived visual angle.

19

# *Relational Graph Construction*

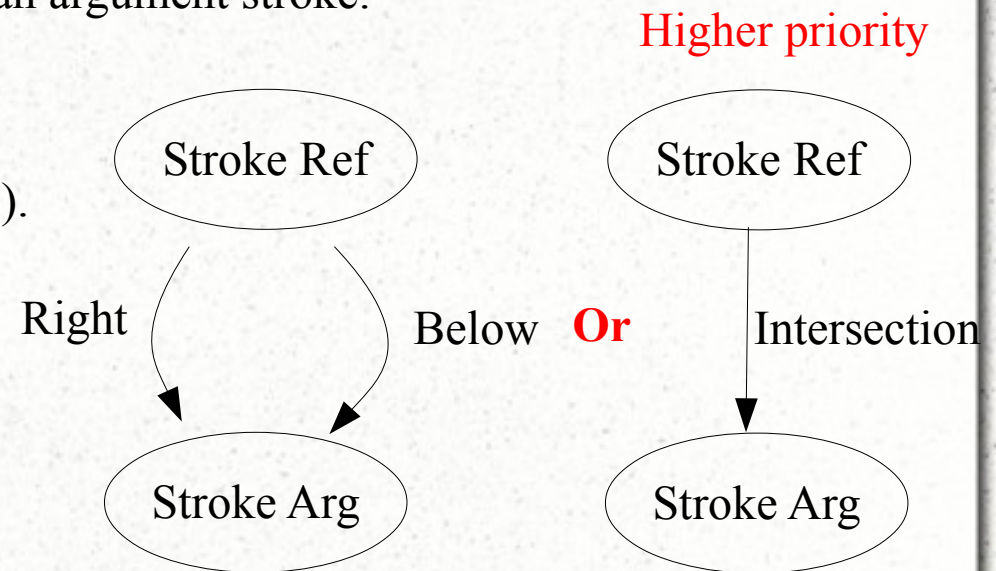Spatial relation: from a reference stroke to an argument stroke.



1. Start with top-left stroke.
2. Choose 2 closest strokes for each spatial relation.
3. Limit relational graph into Directed Acyclic Graph (DAG).
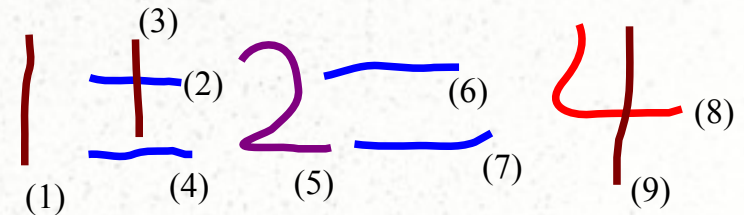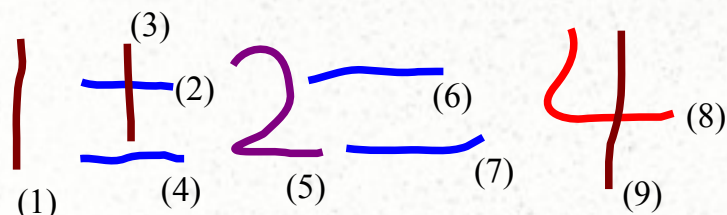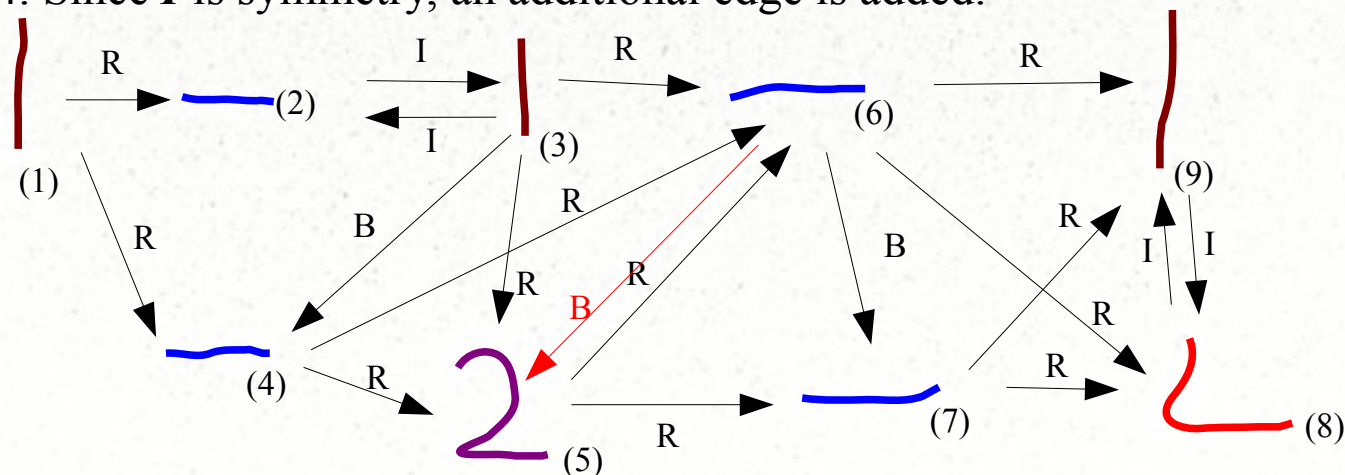4. Since *I* is symmetry, an additional edge is added.

# *Number of edges*

Spatial relation: from a reference stroke to an argument stroke.

$n_r$ : the number of different spatial relations.

$n_{str}$ : the number of strokes.

$n_r = 2 \qquad n_{str} = 9$

Complete directed graph

Number of edges in graph

$$n_r n_{str}(n_{str} - 1)$$ =144

$n_c \le (n_{str} - 1)$

We prefer some symbols composed of the $\boxed{n_c = 2}$ closest strokes

Reduced

$$n_r \cdot n_{str} \cdot n_c$$

$2 \cdot 9 \cdot 2 = 36$

Pruning

21

18

# *Outline*

- 1. Background
- 2. Unsupervised Handwritten Graphical Symbol Learning
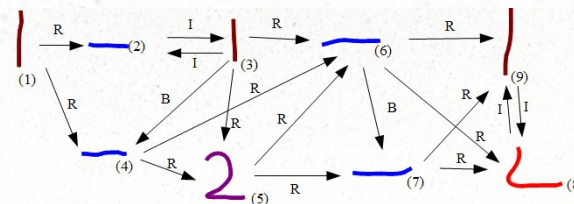    - Quantization of strokes
    - Relational Graph Construction Between Strokes
    - Discover Symbols (Sub-graphs)
- 3. Experiment
- 4. Conclusion

# *Discover Symbols (Sub-graphs)*

Many samples

$4 + 1 = 2$

$4 - 2 = 1$

$1 + 2 = 4$

$4 - 1 = 2$

$2 + 4 = 1$

One sample →

$1 \pm 2 = 4$

(1) (2) (3) (4) (5) (6) (7) (8) (9)

# *Discover Symbols (Sub-graphs)*

Many equations

4 + 1 = 2

4 − 2 = 1

1 + 2 = 4

4 − 1 = 2

2 + 4 = 1

1. We prefer the frequent patterns as the symbols.

=

2. Almost equally frequent pattern
   but with different numbers of strokes.

+    +

Which one?

Three times    Three times

# *Minimum Description Length principle*

Minimum Description Length (MDL) principle is involved in searching the lexical unit that leads to the best compression of data.

[1] describes an unsupervised language learning method using MDL principle on text corpora.

*SUBDUE (SUBstructure Discovery Using Examples) uses the MDL principle to identify patterns that minimize the number of bits needed to describe the input graph after being compressed by the pattern.[2]*

[1] Marcken, C. D., Unsupervised Language Acquisition, Massachusetts Institute of Technology, 1996

[2] Diane J. Cook and Lawrence B. Holder, http://ailab.wsu.edu/subdue/

# *Discover Symbols (Sub-graphs)*

Many samples

4 + 1 = 2
4 - 2 = 1
1 ± 2 = 4
4 - 1 = 2
2 ± 4 = 1

1. We prefer the frequent patterns as the symbols.

2. Almost equally frequent pattern but with different numbers of strokes.

Which one?

Minimum Description Length (MDL) principle

26

# Hierarchical structure (Iterative learning)
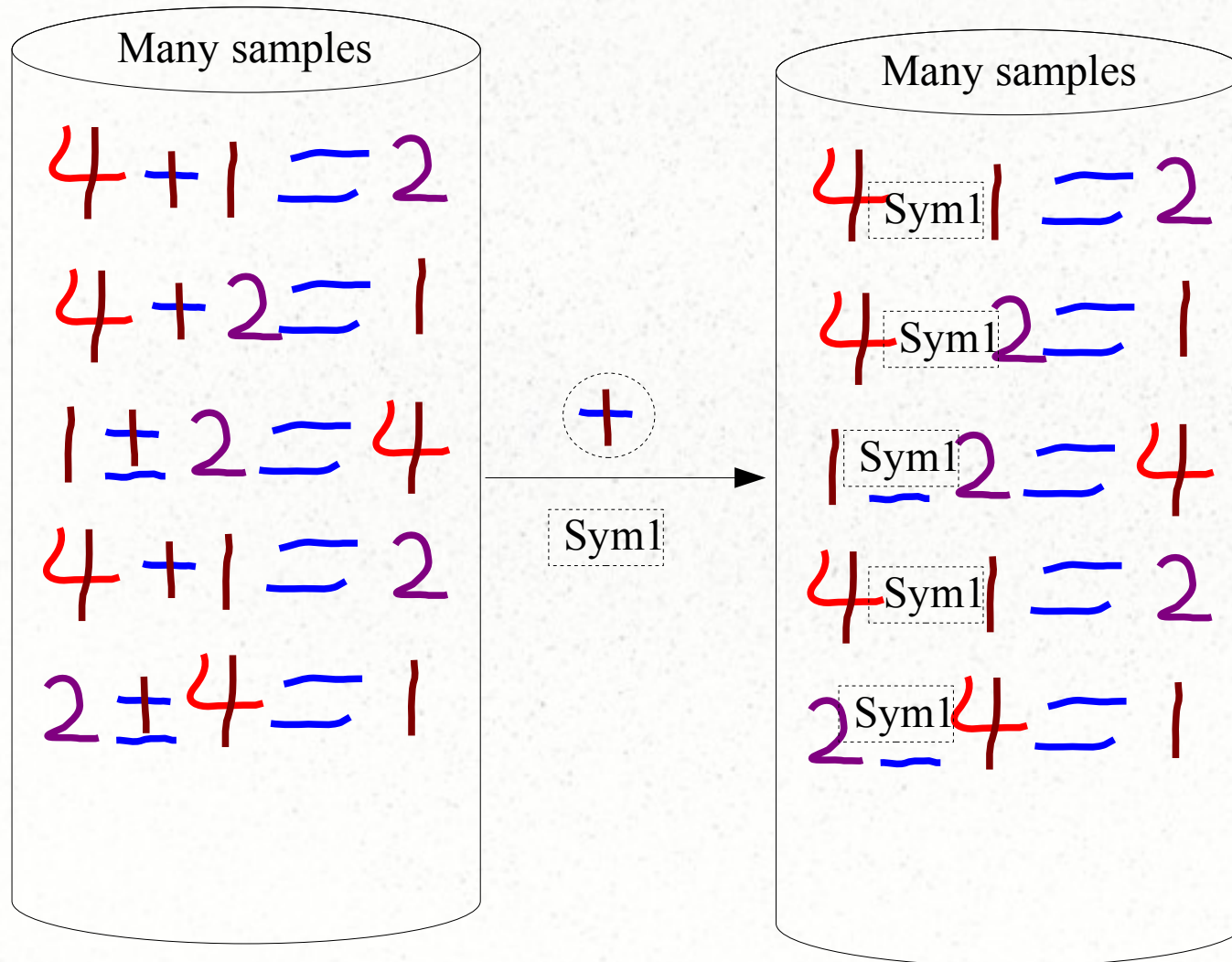
Many samples

$4 + 1 = 2$
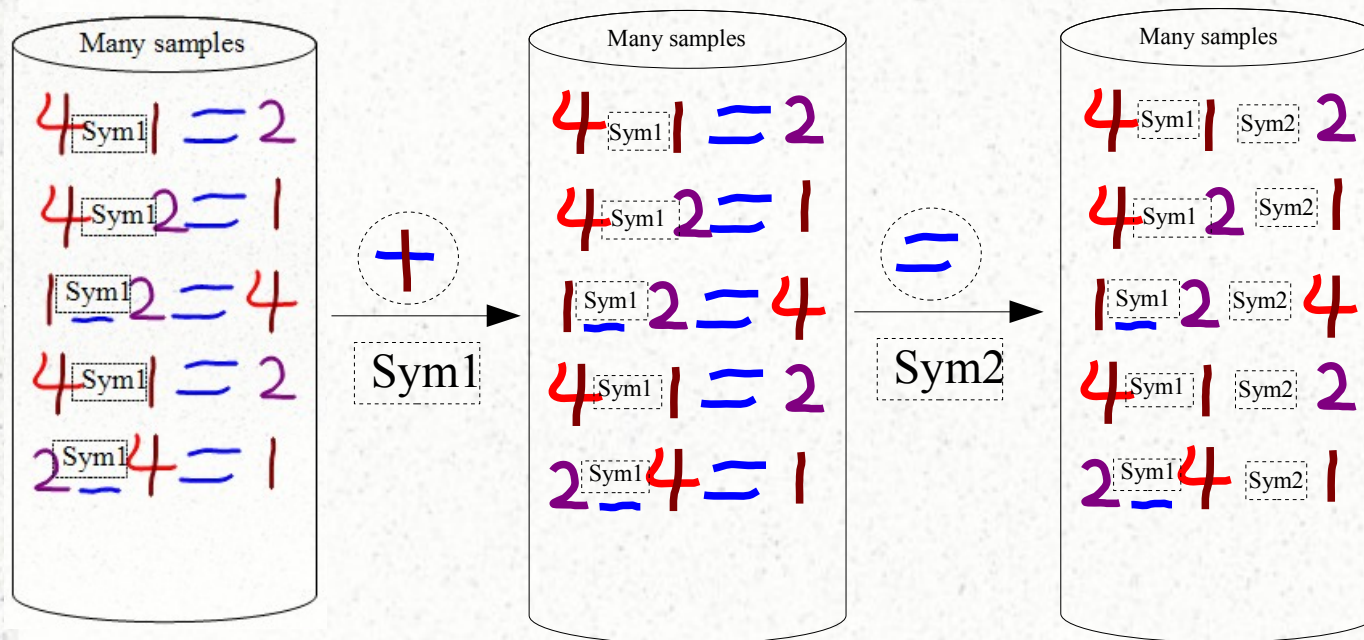
$4 + 2 = 1$

$1 \pm 2 = 4$

$4 + 1 = 2$

$2 \pm 4 = 1$

But if ┿ is much more frequent than ╪ ,

we will choose the ┿ as the symbol according to MDL principle.

# Hierarchical structure (Iterative learning)



Many samples

4 + 1 = 2
4 + 2 = 1
1 + 2 = 4
4 + 1 = 2
2 + 4 = 1

+

Sym1

Many samples

4 Sym1 1 = 2
4 Sym1 2 = 1
1 Sym1 2 = 4
4 Sym1 1 = 2
2 Sym1 4 = 1

28

# Hierarchical structure (Iterative learning)



Lexicon:



Sym1        Sym2

# Hierarchical structure (Iterative learning)



Lexicon:

Sym1    Sym2    Sym3    Sym4

# *Outline*

- 1. Background
- 2. Unsupervised Handwritten Graphical Symbol Learning
  - Quantization of strokes
  - Relational Graph Construction Between Strokes
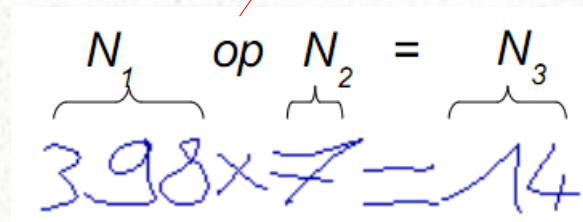  - Discover Symbols (Sub-graphs)
- 3. Experiment
- 4. Conclusion

# *Dataset (Experiment)*

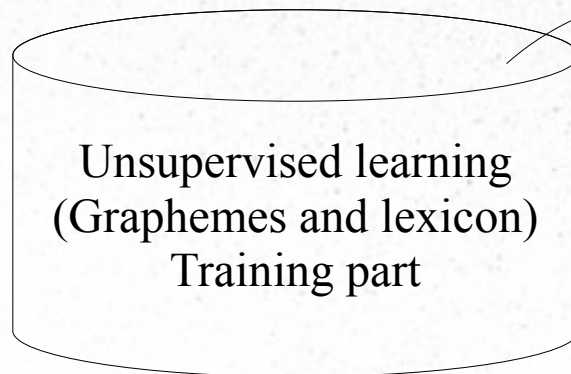A synthetic dataset from real isolated handwritten characters.

$N_{i=\{1,2,3\}}$ is 70% of 1 digit, 20% of 2 digits and 10% of 3 digits randomly.
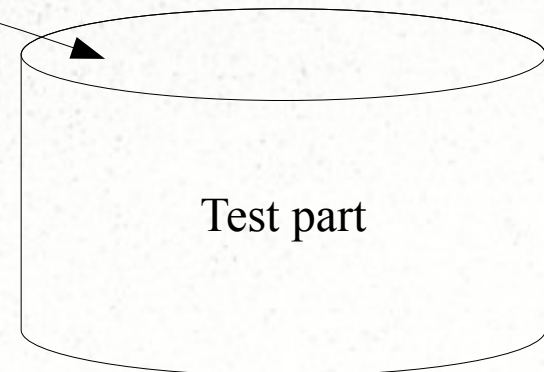
$\{0, 1, ..., 9\}$ $\{+, -, \times, \div\}$



$$N_1 \quad op \quad N_2 \quad = \quad N_3$$

Lexicon

Unsupervised learning
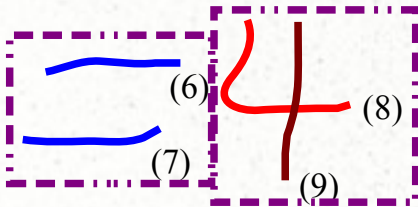(Graphemes and lexicon)
Training part

Test part

5427 symbols from 180 writers
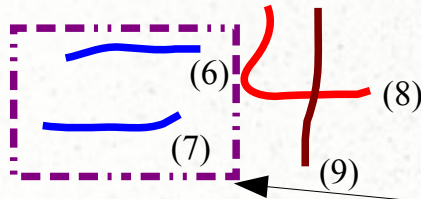
3035 symbols from 100 writers

# *Measure (Experiment)*



*S(e, G)*:ground-truth for the expression.

$$S(e, G) = \{\{(6),(7)\}, \{(8),(9)\}\}$$

# *Measure (Experiment)*



$S(e, G)$:ground-truth for the expression.

$$S(e,G) = \{\{(6),(7)\},\{(8),(9)\}\}$$

$S(e, L)$:hierarchical segmentation using lexicon $L$.

$$S(e,L) = \{\{(6)\},\{(7)\},\{(6),(7)\},\{(8)\},\{(9)\}\}$$

$$R_{\text{Recall}} = \frac{|S(e,G) \cap S(e,L)|}{|S(e,G)|} = 0.5$$

We got the recall rate of **64.3%** for **multi-stroke** symbols
(863 symbols from 1343 symbols) on the test part of our dataset.

34

# *Conclusion*

- Quantization of strokes

- Construction of relational graph

- Lexicon extraction using MDL principle (SUBDUE)

- The recall rate of 64.3% (863/1343 multi-stroke symbols) is obtained.

# *Future work*

- More complex spatial relation definition for more complex language, such as flowchart.

- Annotation assistance system for graphical symbols

Thank you for your attention.
Questions?