

A. Pastor López-Monroy<sup>1,2</sup>, Fabio A. González<sup>3</sup>

Manuel Montes-y-Gómez<sup>4</sup>, Hugo Jair Escalante<sup>4</sup>, Thamar Solorio<sup>2</sup>

<sup>1</sup>Mathematics Research Center CIMAT, <sup>2</sup>Department of Computer Science, University of Houston, USA

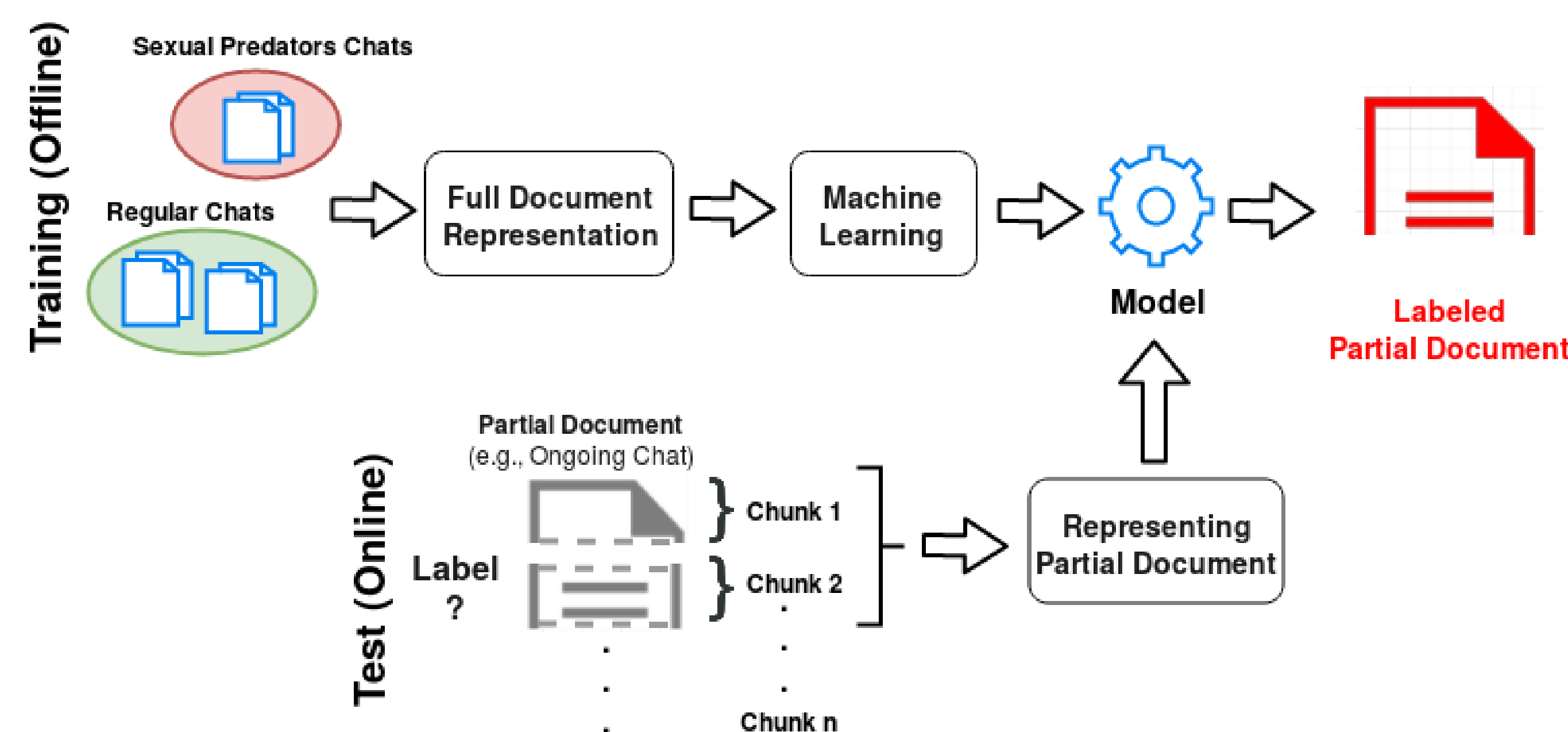
<sup>3</sup>Systems and Computer Engineering Department, Universidad Nacional de Colombia

<sup>4</sup>Department of Computer Science, INAOE, México

## 1. Introduction

- Detecting risks on Social Media data with as much anticipation as possible is crucial for prevention.
- **Early Text Classification (ETC)** is an emerging field where the goal is to anticipate the prediction by using as little text as possible.

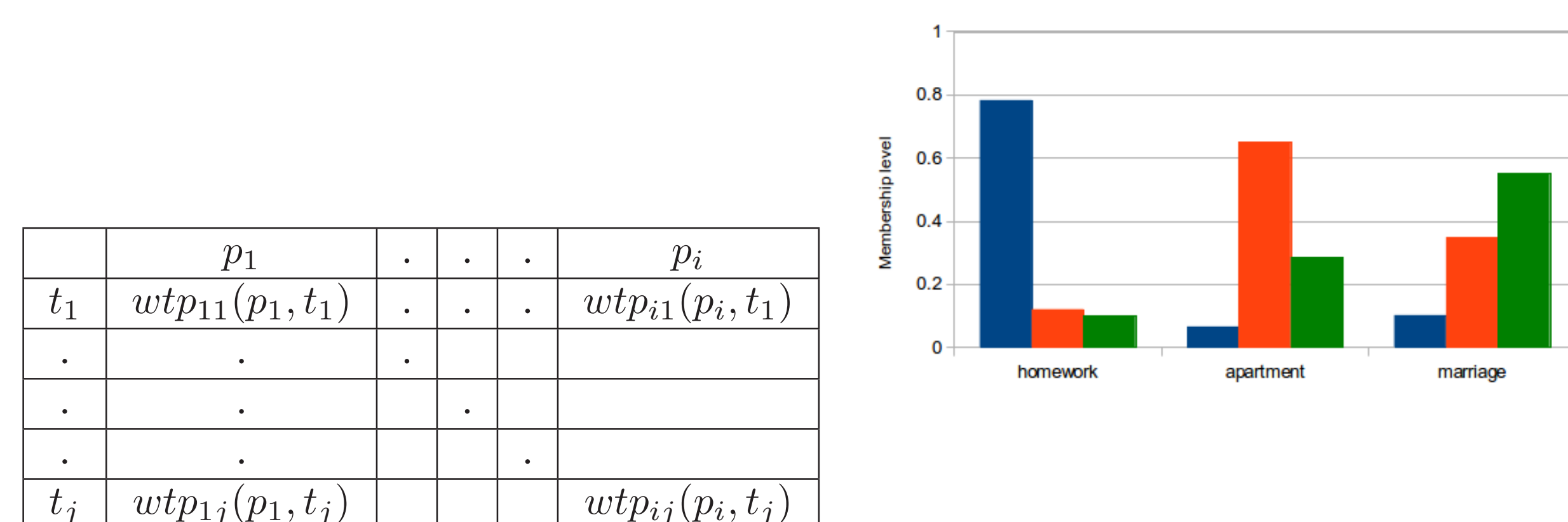
### Early Classification on Social Media



- We need to model very short length documents on early stages.
- Late stages require to exploit the additional evidence as much as possible to make accurate predictions.

## 2. Word Representation

### Temporal Variation of Terms (TVT)



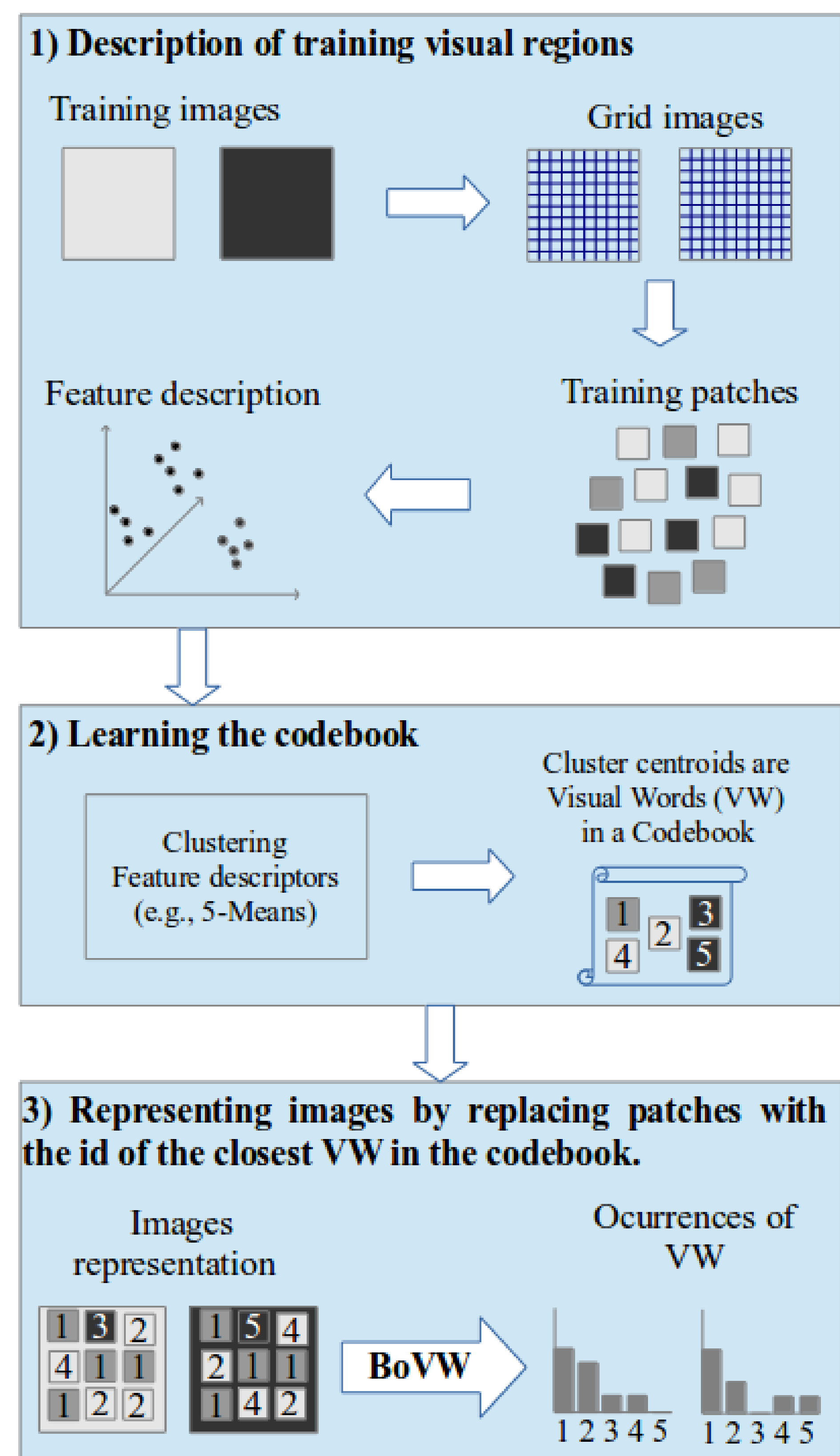
The representation builds term vectors  $\mathbf{t}_j = \langle tp_{1j}, \dots, tp_{ij} \rangle$ , where  $tp_{ij}$  is a value representing the relationship of the term  $t_j$  with the class  $p_i$ .

$$wt_{p_{ij}} = \sum_{k: d_k \in P_i} \log_2 \left( 1 + \frac{tf_{kj}}{\text{len}(d_k)} \right)$$

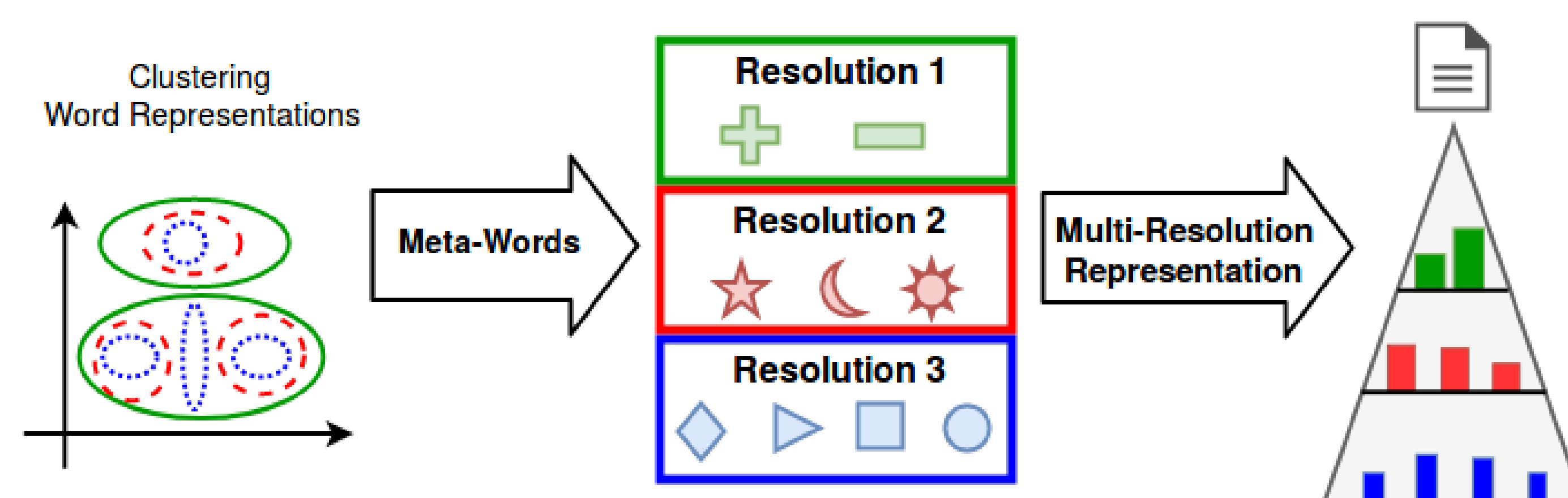
These word vectors are built at specific chunks in training.

## 3. Multi Resolution Representation (MulR)

Single Resolution is analogous to the BoVW

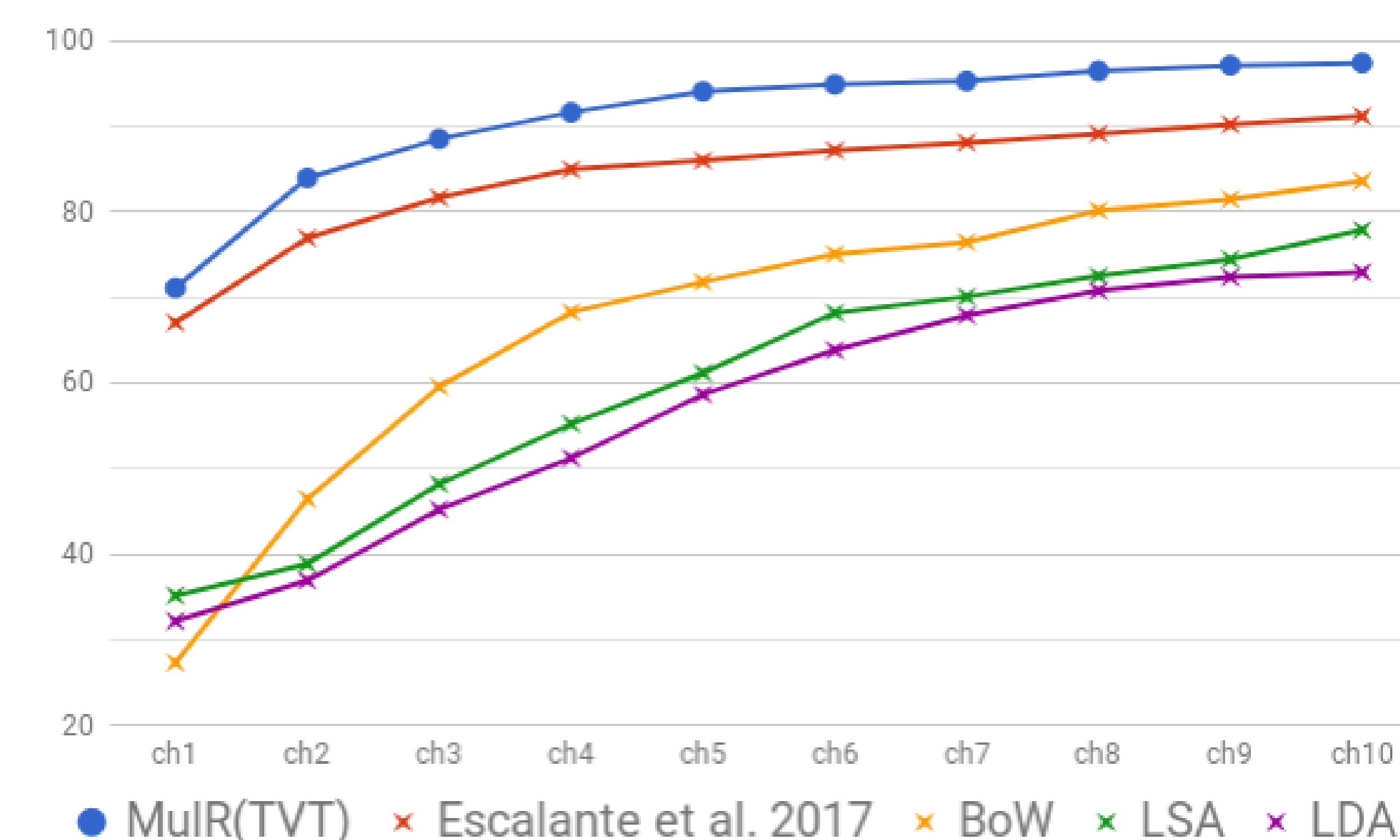


Intuitive idea of Multi-Resolutions



## 4. Evaluation

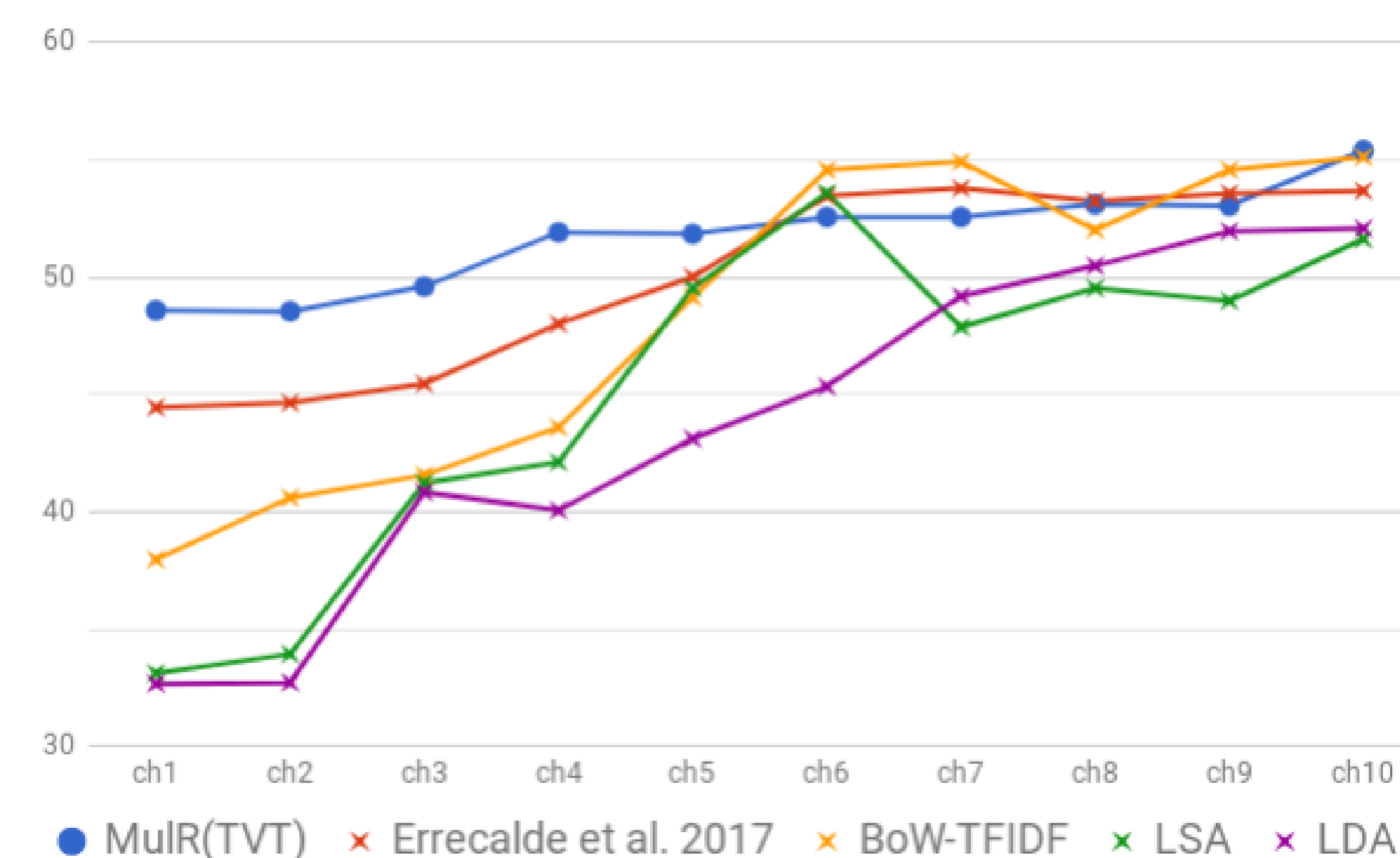
### Sexual Predator Detection



1. The best word vector representation was TVT for ETC.
2. MulR is effective for all early chunks.
3. As more text is available, methodologies significantly improve.

Test set	$R_1$	$R_2$	$R_3$	$R_4$	$R_5$
chunk-1	4	3	2	1	0
chunk-2	3	2	3	2	0
chunk-3	3	2	2	2	1
chunk-4	3	2	3	1	1
chunk-5	3	2	3	1	1
chunk-6	3	1	4	1	1
chunk-7	3	1	3	2	1
chunk-8	2	2	1	2	3
chunk-9	3	1	1	2	3
chunk-10	3	0	2	2	3

### Depression Detection



1. MulR(TVT) improves between  $\approx 5\%$  and  $\approx 2\%$  in chunks 1 to 4.
2. Depression Detection problem is much harder than SPD.
3.  $F_1$  under  $\approx 60\%$  and highly unbalanced dataset.