

## ON THE APPLICABILITY OF BINARY CLASSIFICATION TO DETECT MEMORY ACCESS ATTACKS IN IOT

C&ESAR 2018- Rennes | CEA Leti | KERROUMI Sanaa | 08/11/18

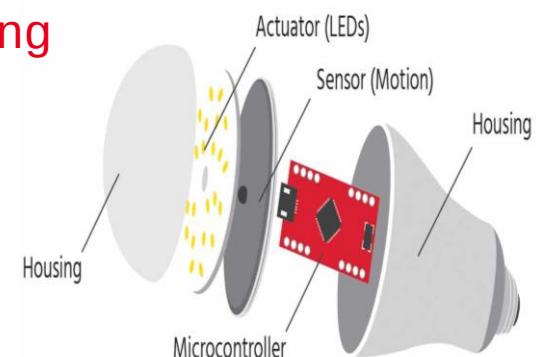
- 1** IoT node
- 2** Related works
- 3** Problem statement
- 4** Proposed methodology
- 5** Results
- 6** Take out and lessons learned

- Internet Of Things

“The interconnection via the internet of computing devices embedded in everyday objects enabling them to communicate”

- The “thing” in IoT can be anything and everything as long as it has a unique identity and can communicate via the internet

- Sensors, actuators or combined sensor/actuator
- Limited capabilities in terms of their computational power, memory, energy, availability, processing time, cost, ... → limits their abilities to handle encryption or other data security functions
- Designed to be disposable → updates/security patches may be difficult or impossible.
- Designed to last for decades → any unpatched vulnerabilities will stay for very long
- A foothold in the network (e.g., IoT goes nuclear, thermometer in the fish tank attack)



# EDGE-NODE VULNERABILITIES: WHAT COULD POSSIBLY GO WRONG?

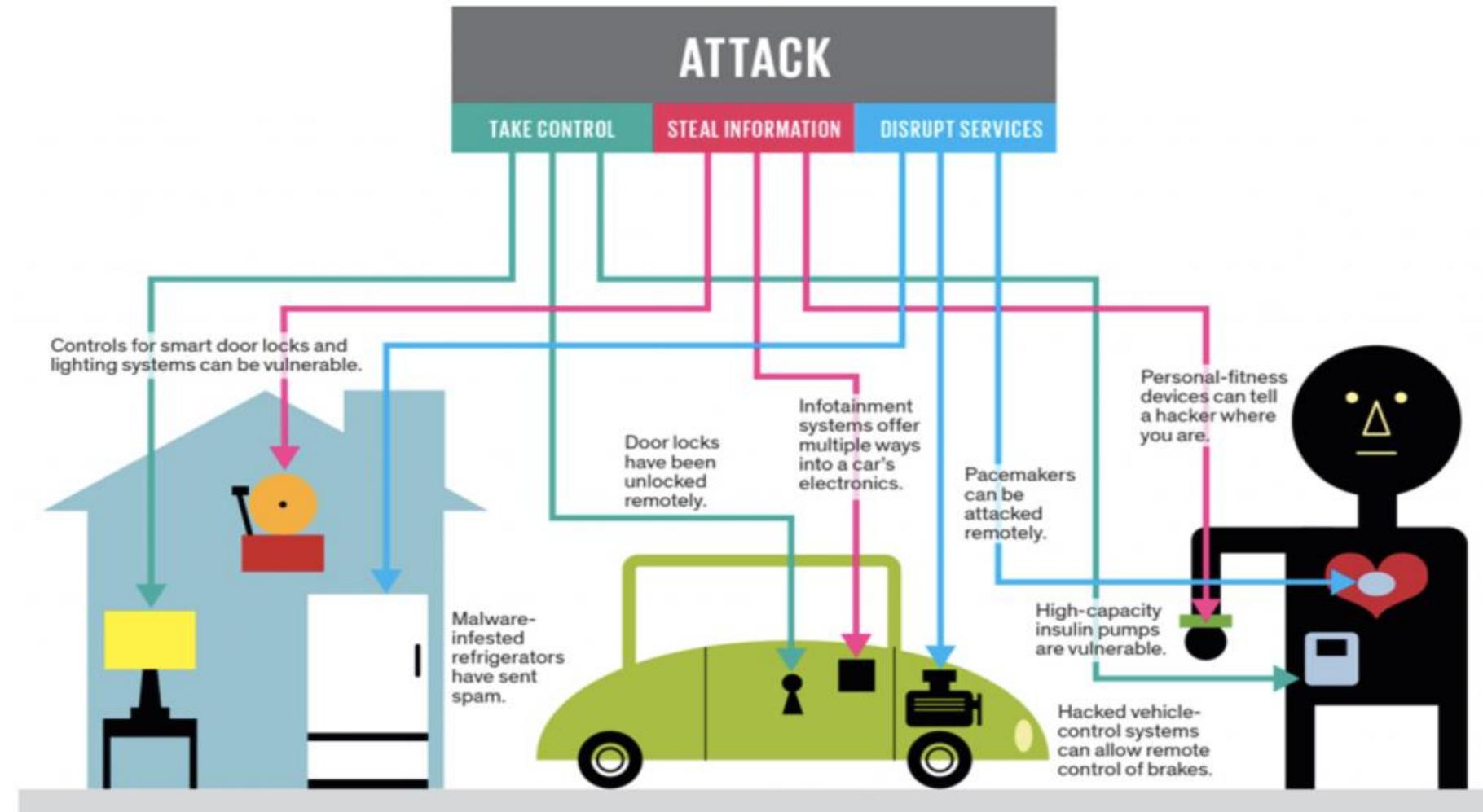
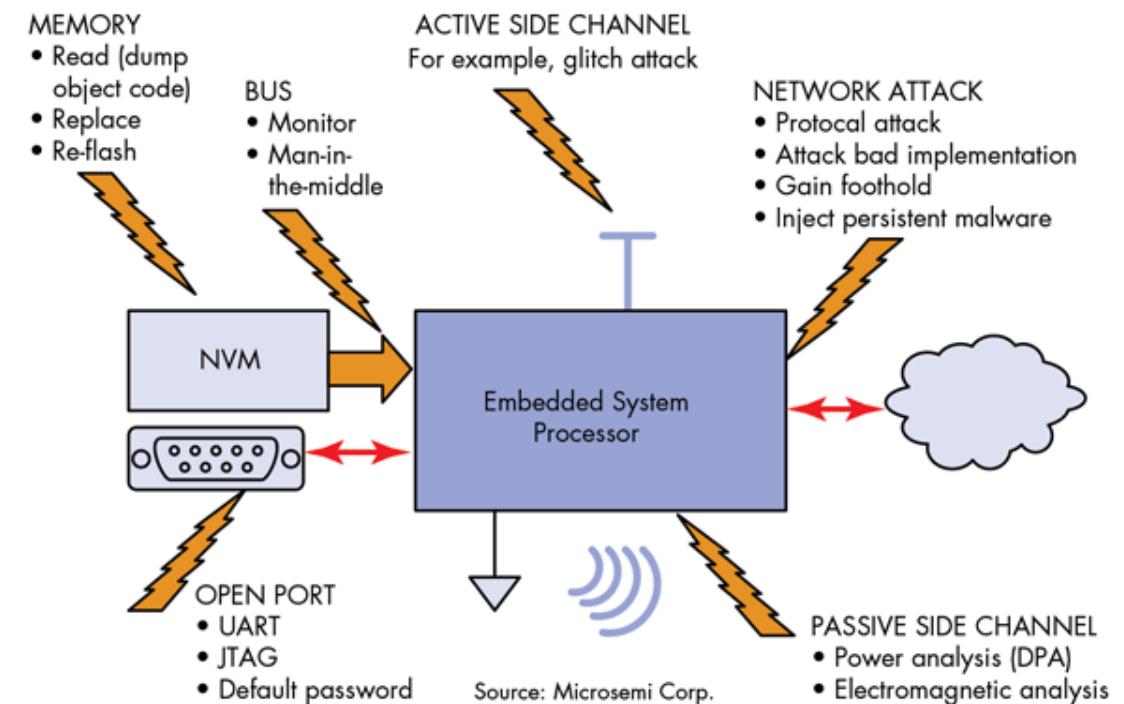


Illustration: J. D. King

- Attack modes:
  - Software attacks
  - Side channel attacks
  - Physical attacks
  - Network attacks
- Why are we interested in the memory access attacks?
  - It is particularly hard to fake or hide malicious tasks memory accesses
  - It offers a great view on what's going on inside the device
  - Alluring target for the attacker
    - Control the node
    - Read encryption keys or protected code
    - ...



Fuses and flash readout protection	Encryption	Detection
<ul style="list-style-type: none"><li>• <u>Pros</u><ul style="list-style-type: none"><li>• Inexpensive</li><li>• Efficient</li><li>• Easy to implement</li></ul></li><li>• <u>Cons</u><ul style="list-style-type: none"><li>• Mostly set on level that permits access to memory (post deployment upgrades)</li></ul></li></ul>	<ul style="list-style-type: none"><li>• <u>Pros</u><ul style="list-style-type: none"><li>• Preserve privacy and confidentiality</li><li>• Convenient</li></ul></li><li>• <u>Cons</u><ul style="list-style-type: none"><li>• Expenses</li><li>• Compatibility</li><li>• Encryption keys</li><li>• Widespread security compromise</li></ul></li></ul>	<ul style="list-style-type: none"><li>• <u>Pros</u><ul style="list-style-type: none"><li>• Proactive</li><li>• Scalable</li><li>• Pervasive</li><li>• Great 1<sup>st</sup> line of defense</li></ul></li><li>• <u>Cons</u><ul style="list-style-type: none"><li>• False Positives</li><li>• Leak</li><li>• Mimicry attacks</li></ul></li></ul>

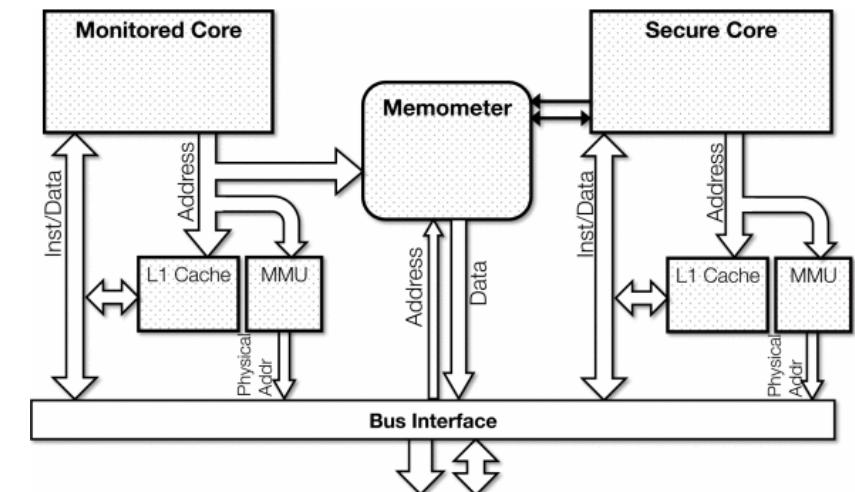
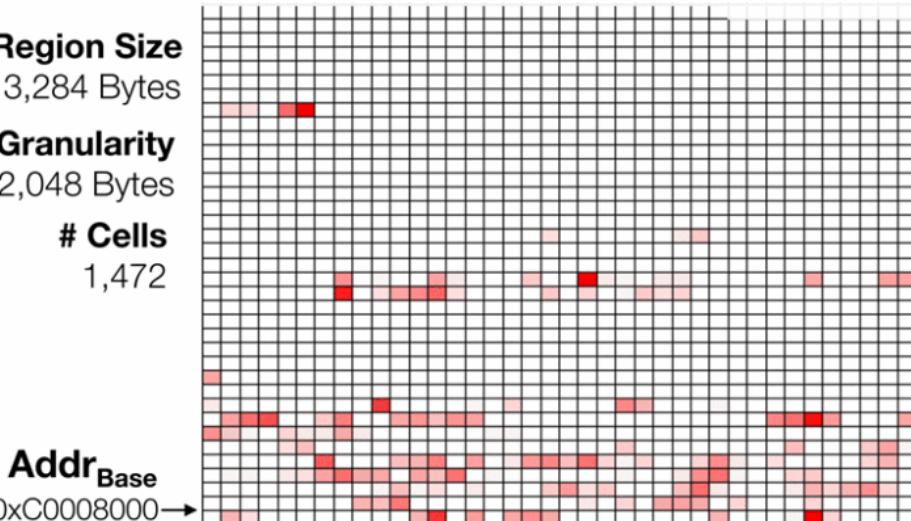
- **Memory heat map**

- Idea: profiling memory behavior by representing the frequency of access to a particular memory region (regardless of which component accessed it) during a time interval. The MHM is then combined with an image recognition algorithm to detect any anomalies.
- Strengths:
  - system wide anomalies detection (not just malicious ones)
  - can be used in real-time embedded systems
- Limitations :
  - expensive to compute: need to store several images of nominal MHM
  - wrong architecture (Config3 and higher )

**Memory Region Size**  
3,013,284 Bytes

**Granularity**  
2,048 Bytes  
**# Cells**  
1,472

**Addr Base**  
0xC0008000→



Yoon, Man Ki, et al, "Memory heat map: anomaly detection in real-time embedded systems using memory behavior". In *Design Automation Conference (DAC)*, 2015 52nd ACM/EDAC/IEEE (pp. 1-6). IEEE.

## System call distribution

- Idea : learn the normal system call frequency distributions, collected during legitimate executions of a sanitized system, combined by a clustering algorithm (k-means). If an observation, at a run time, is not similar to any identified clustered. The observation is doomed malicious
- Strengths:
  - simple
- Limitations
  - Require an OS
  - need a throughout training
  - no adaptation of centroids (any change even if nominal would be flagged as malicious)
  - application to be monitored need to be very deterministic
  - definition of cut off line influence the FPR and detection rate

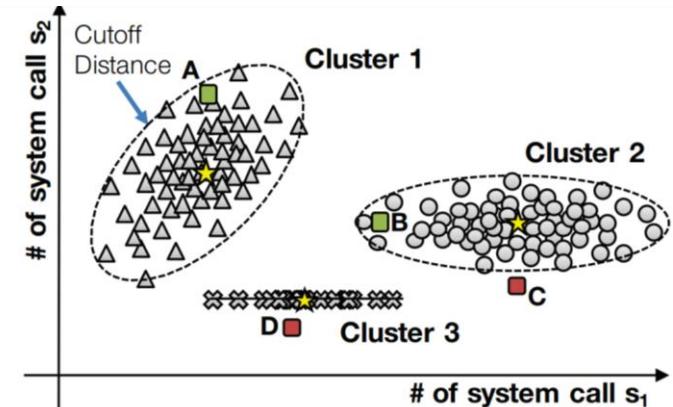
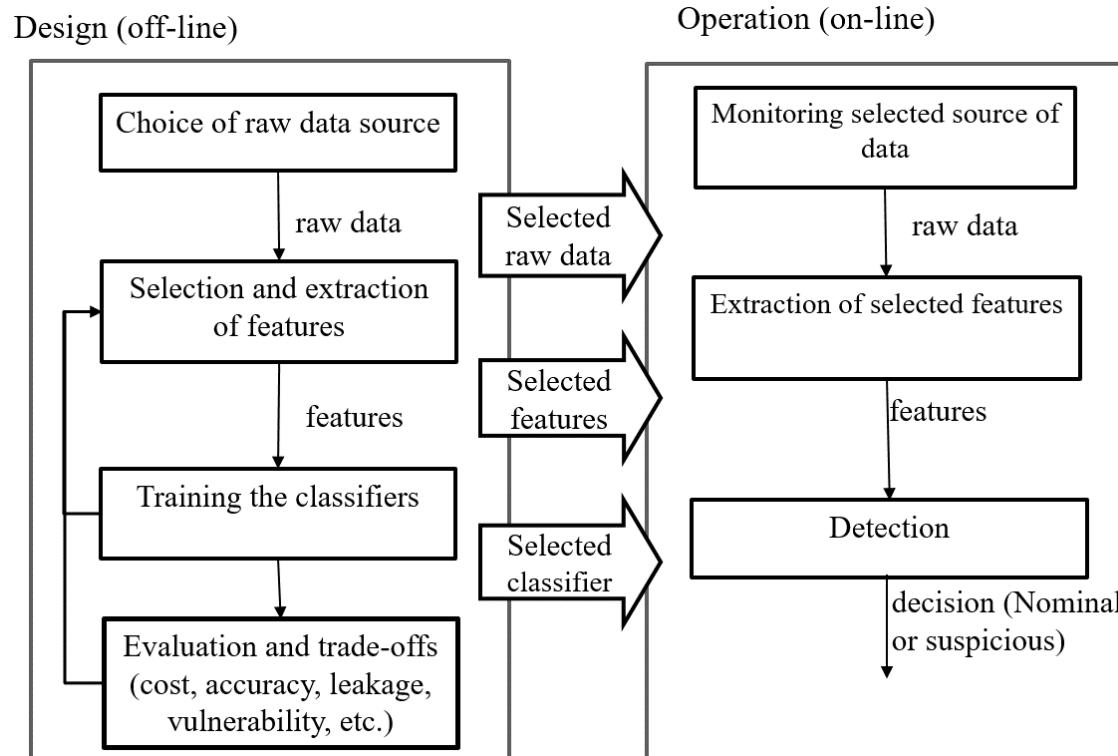


Figure 3: System call frequency distributions for  $S = \{s_1, s_2\}$  and clusters. The gray-colored objects are SCFDs in the training set. Each star-shaped point is the centroid of each cluster. The ellipsoid around each cluster draws its cutoff line.

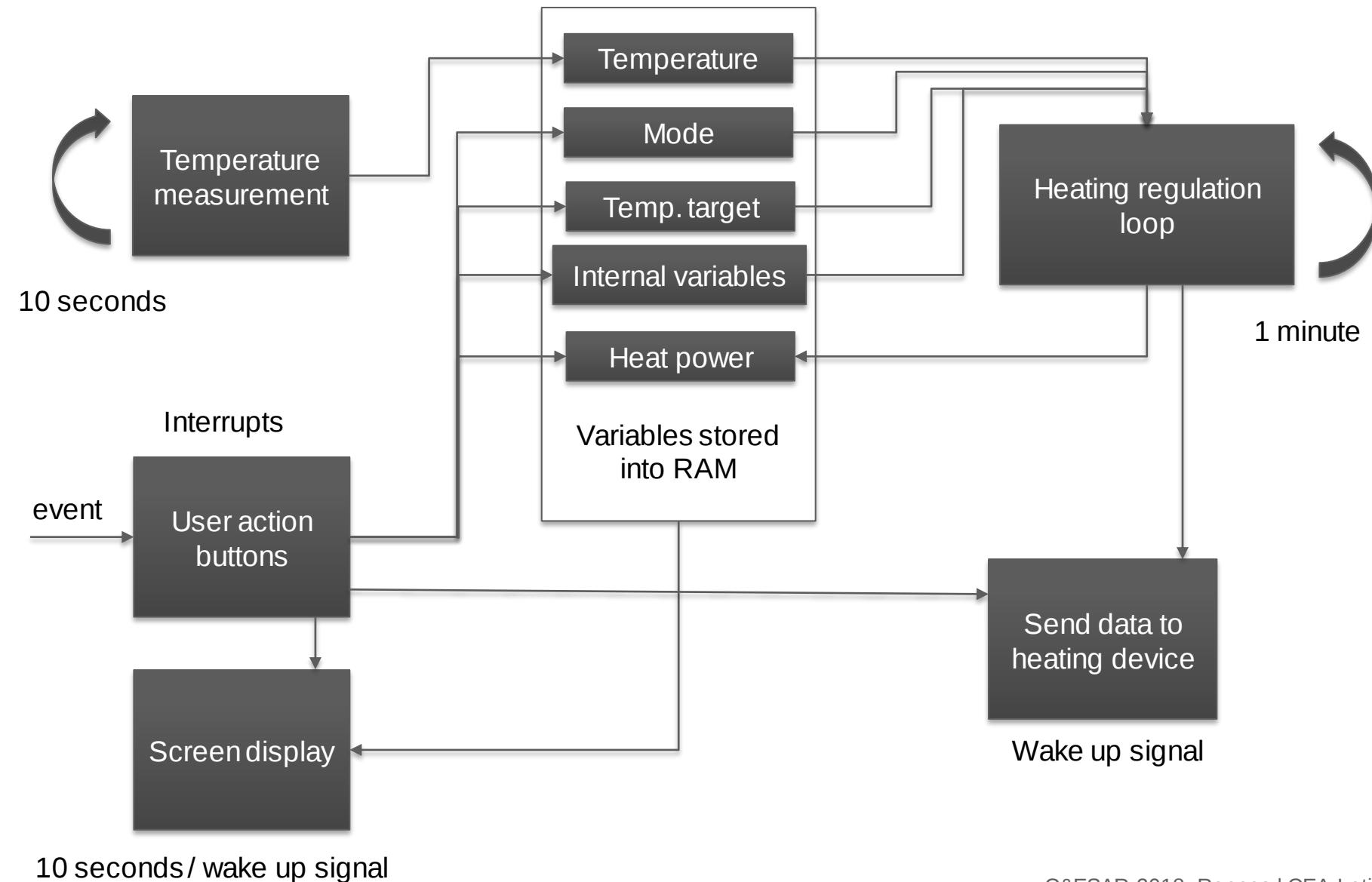
- **Existent detection solutions are:**
  - Not directly related to memory access attacks
  - Too expensive to compute
  - Used features in detection are either hard or impossible to acquire for constrained node (e.g., hardware performance counters, control flow, instruction mix, etc.)
- **Analyze the effectiveness of binary classifiers combined by simples features to detect memory access attacks in the context of a low cost IoT node**

- 2 phases' methodology:
  - Design: performed during the design of the node to build the detector
  - Operation: the detector in operation



- In this presentation we will focus on the design part of the detector

## USE CASE PRESENTATION: CONNECTED THERMOSTAT



## IN MORE DETAILS

Processor/  
memory  
trace

Feature  
extraction  
& selection

Machine  
learning  
method

Evaluation  
and trade-  
offs

- Raw data: memory access log
  - Timestamp
  - Accessed address
  - Data manipulated
  - Type of data
  - Flag to indicate if the access is nominal or suspicious
- Features – computed each time window
  - Number of memory reads, number of memory accesses, cycles between consecutive reads, address increment, number of “unknown” (first-encountered) addresses, amount of read/accessed data ...

0000000001 READ32	00000018	->	00000003
0000000001 WRITE8	000008f1	<-	00
0000000001 READ32	00000013	->	00000003
0000000001 WRITE8	000008f1	<-	00
0000000001 READ32	00000000	->	00000000
0000000015 READ32	00000004	->	00000000
0000000041 READ32	00000008	->	00000000
0000000065 READ32	0000000c	->	00000000
0000000088 READ32	00000010	->	4180004001
0000000111 READ32	00000014	->	41980000
0000000134 READ32	00000018	->	00000003
0000000157 READ32	0000001c	->	00000000
0000000180 READ32	00000020	->	00000000
0000000203 READ32	00000024	->	00000000
0000000226 READ32	00000028	->	00000000
0000000249 READ32	0000002c	->	00000000
0000000272 READ32	00000030	->	00000000
0000000295 READ32	00000034	->	00000000
0000000318 READ32	00000038	->	00000000
0000000341 READ32	00000040	->	00000000
0000000364 READ32	00000044	->	00000000
0000000387 READ32	00000048	->	00000000
0000000411 READ32	0000004c	->	00000000
0000000434 READ32	00000050	->	00000000
0000000457 READ32	00000054	->	00000000
0000000474 READ32	00000058	->	00000000
0000000521 READ32	0000045d	->	00000000
0000000544 READ32	0000043e	->	00000000
0000000567 READ32	00000442	->	00000000

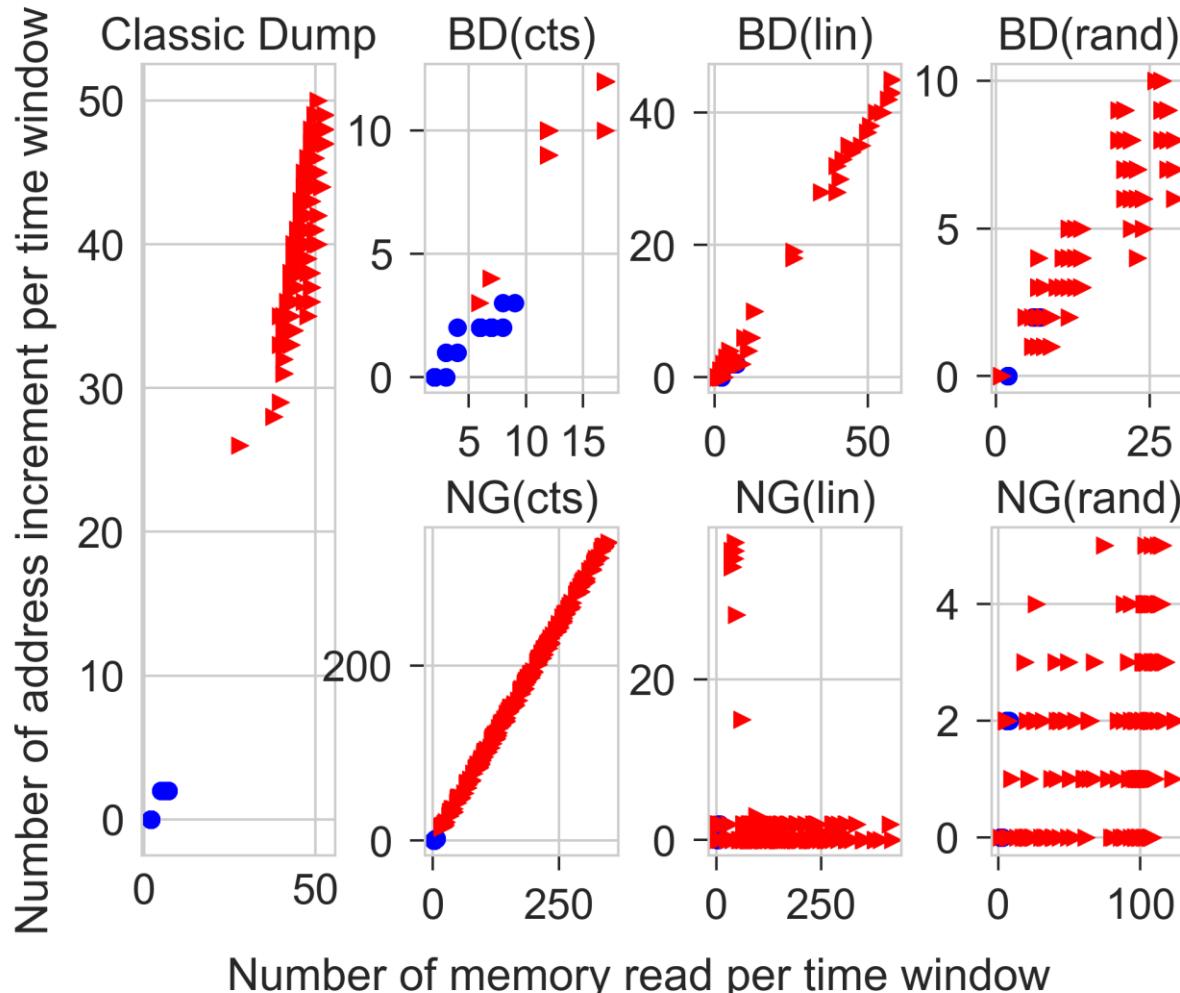
Detected!

- **Classic dump (CD): basic memory dump require minimal effort from the attacker:**
  - Attacker reads the entire memory in a contiguous way, the memory reads are spaced regularly in time and memory space

Attacker assumed to be aware of the presence of some security monitor → avoid obvious change in the memory patterns of the device

- **Dumping in bursts (DB):**
  - The memory is read in bursts, the accessed addresses are still contiguous but the time step between two consecutive reads is incremented by constant (BD(cts)), linearly (BD(lin)) or randomly (BD(rand))
- **Dump in non contiguous way (NG)**
  - The address increment between two consecutive reads is incremented by constant (NG(cts)), linearly (NG(lin)) or randomly (NG(rand))

## TRAINING &amp; TESTING DATASETS



dataset	Training	Testing
Experiment 1	Nominal + CD	DB and NG
Experiment 2	(1) Nom+(CD+NG+BD) (2) Nom+(CD+NG) (3) Nom+(CD+BD)	(1) Nom+(CD+NG+BD)* (2) Nom+BD (3) Nom+ NG

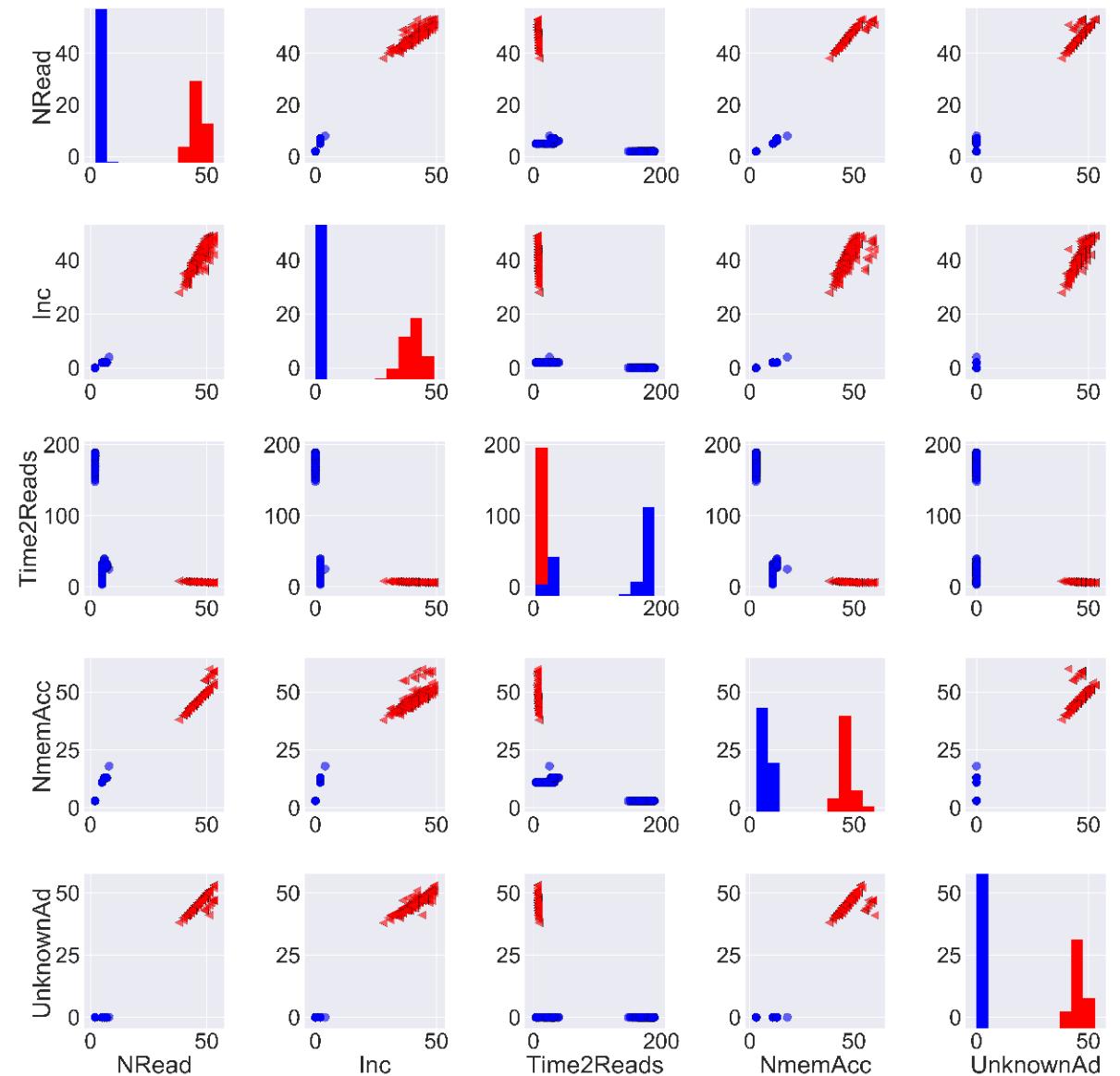
## EXTRACTED FEATURES

Processor/  
memory  
trace

Feature  
extraction  
& selection

Machine  
learning  
method

Evaluation  
and trade-  
offs



Type  
● Nominal  
● Attack

- Nread: number of reads per time interval
- Inc: number of address increment per time interval
- Time2Reads: average time elapsed between two consecutive reads in time interval
- NmemAcc: number of memory access per time interval
- UnknownAd: number of unknown addresses accessed during a time interval

# CLASSIFICATION



Processor/  
memory  
trace

Feature  
extraction  
& selection

Classifiers

Evaluation  
and trade-  
offs

- Let  $X=\{x_1, \dots, x_n\}$  be our dataset and let  $y_i \in \{1, -1\}$  be the class label of  $x_i$
- The decision function ( $f$ ) assign each new instance a label based on prior knowledge gathered during the training
- List of classifiers included in the analysis**
  - K nearest neighbor, Support vector machine, decision tree, random forest, naïve Bayes, linear discriminant analysis and quadratic discriminant analysis

- **Assumption:**
  - Features are independent
- **Intuition:**
  - Given a new unseen instance, we (1) find its probability of it belonging to each class, and (2) pick the most probable.

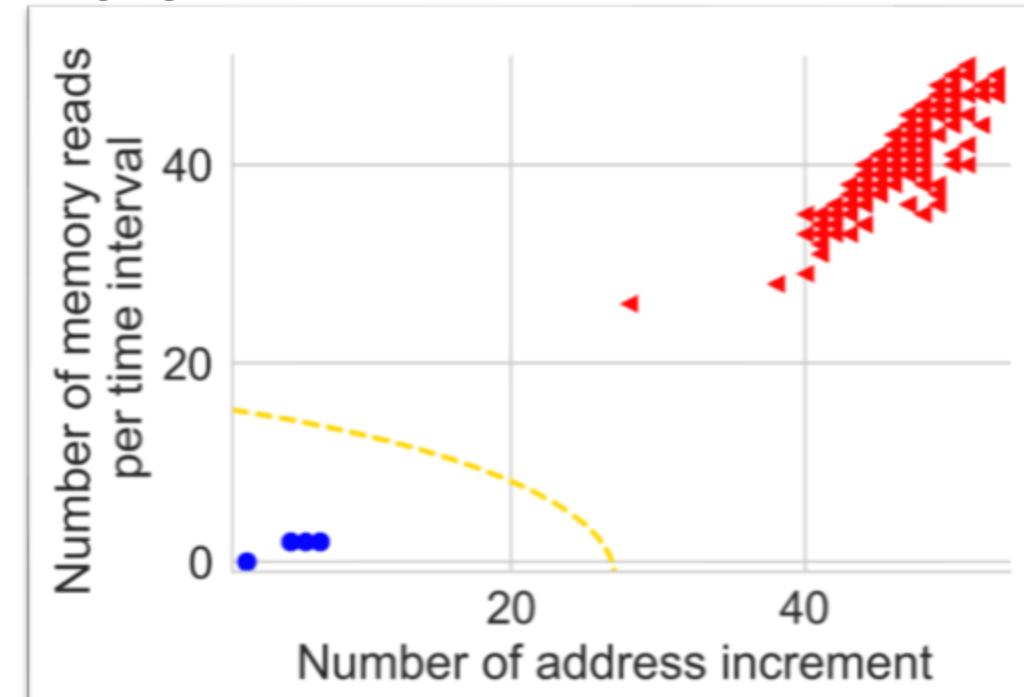
$$P(c_j|x) = \frac{P(x|c_j)P(c_j)}{p(x)}$$

Posterior probability

Likelihood

Predictor prior probability

Class prior probability



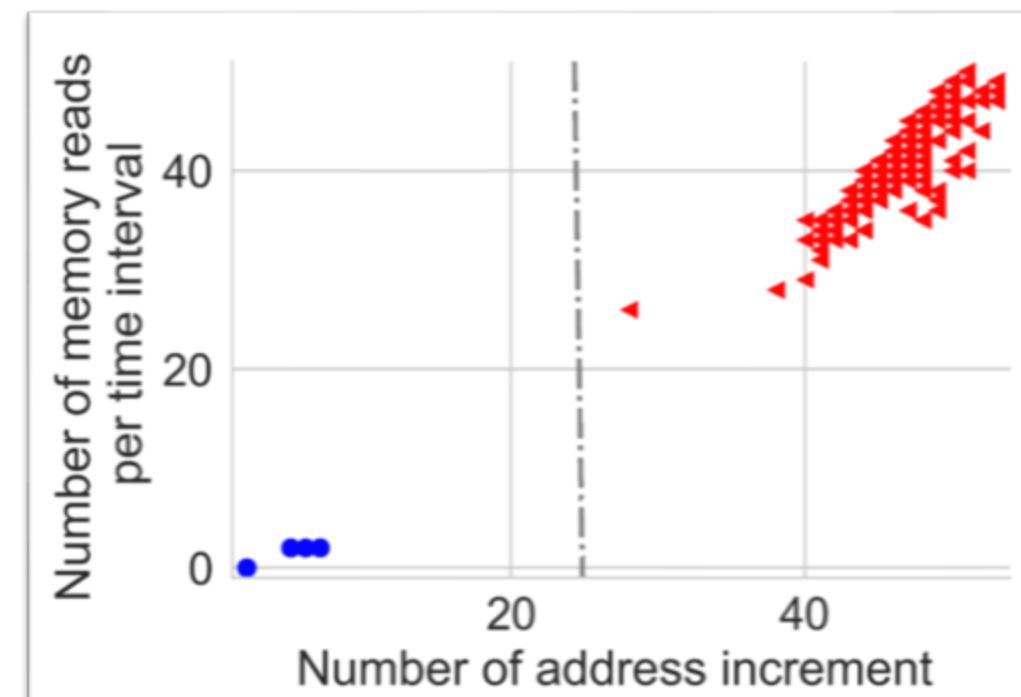
## • Assumption :

- Every class distribution is Gaussian and the covariance matrices are identical

## • Intuition

$$\hat{\delta}_{k(x)} = x^T \Sigma^{-1} \hat{\mu}_k - \frac{1}{2} \hat{\mu}_k^T - \hat{\mu}_k + \log(\hat{\pi}_k)$$

- $\hat{\delta}_{k(x)}$  is the estimated discriminant score that the observation will fall in the  $k^{\text{th}}$  class based on the value of the predictor variable  $x$
  - $\hat{\mu}_k$  is a class-specific mean vector, and  $\Sigma$  is a covariance matrix that is common to all  $K$  classes
  - $\hat{\pi}_k$  is the prior probability that an observation belongs to the  $k^{\text{th}}$  class
- An observation will be assigned to class  $k$  where the discriminant score  $\hat{\delta}_{k(x)}$  is the largest,



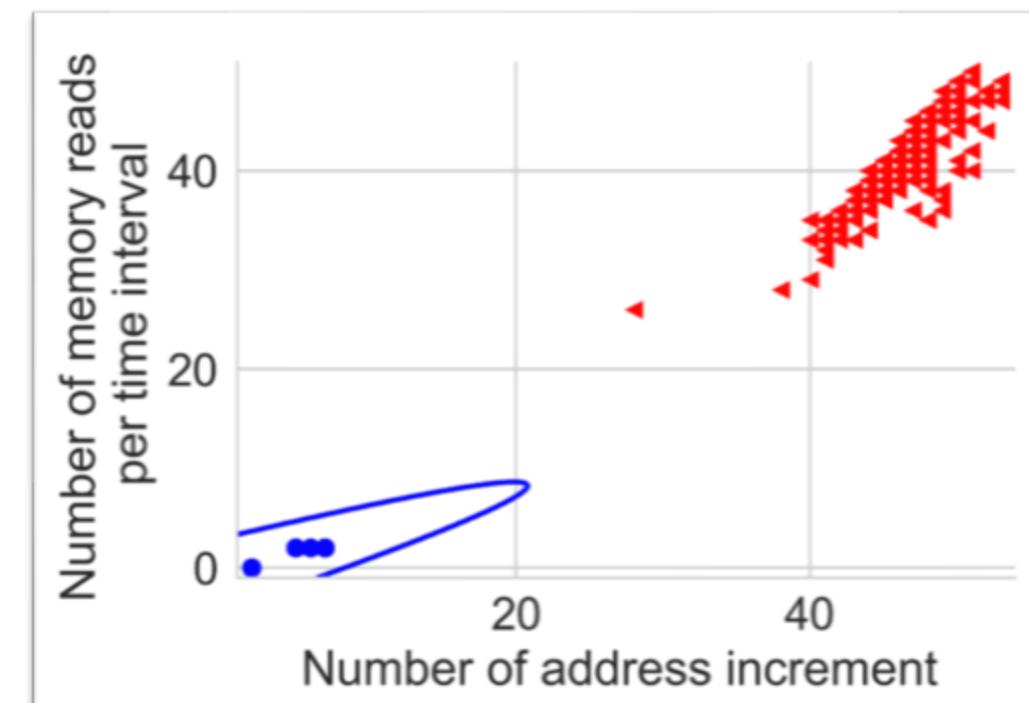
- Assumptions

- Class distribution is Gaussian but with different covariance matrices

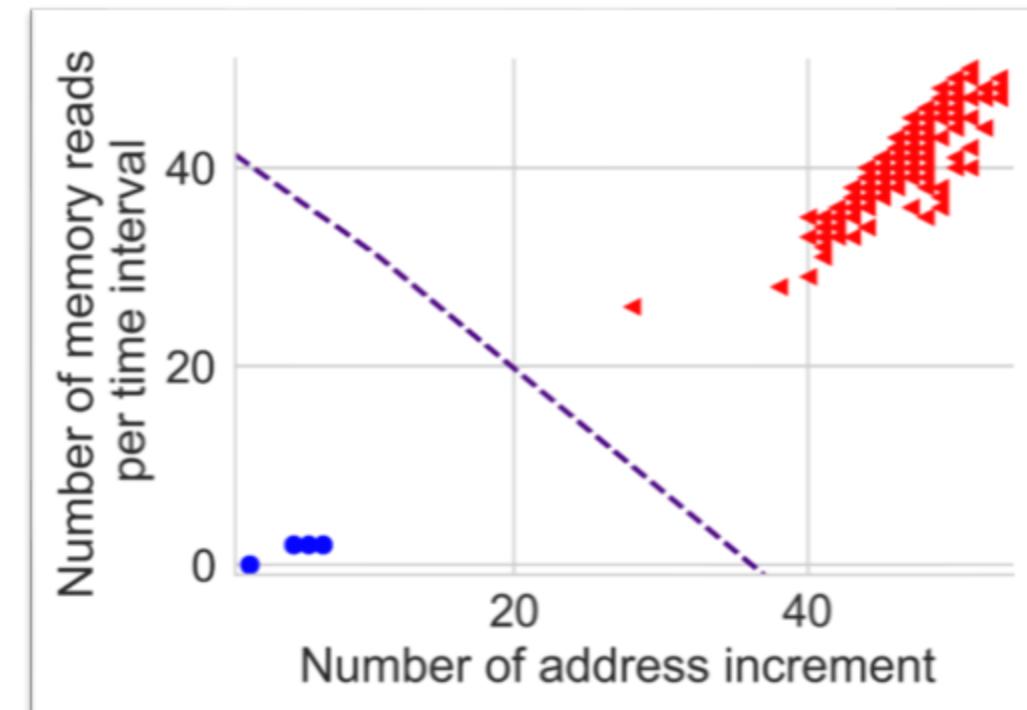
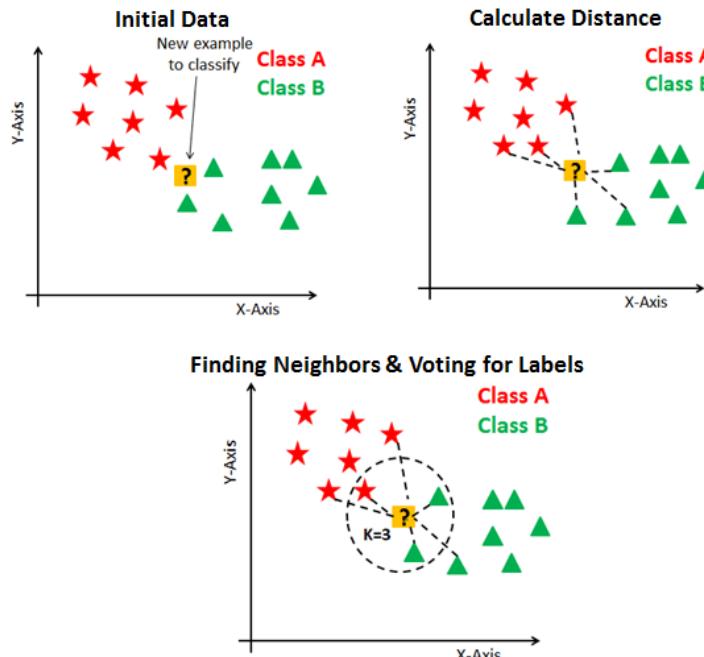
- Intuition

$$\hat{\delta}_{k(x)} = x^T \Sigma^{-1} \hat{\mu}_k - \frac{1}{2} \hat{\mu}_k^T - \hat{\mu}_k + \log(\hat{\pi}_k)$$

- $\hat{\delta}_{k(x)}$  is the estimated discriminant score that the observation will fall in the  $k^{\text{th}}$  class based on the value of the predictor variable  $x$
- $\hat{\mu}_k$  is a class-specific mean vector, and  $\Sigma$  is a covariance matrix that is common to all  $K$  classes
- $\hat{\pi}_k$  is the prior probability that an observation belongs to the  $k^{\text{th}}$  class
- an observation will be assigned to class  $k$  where the discriminant score  $\hat{\delta}_{k(x)}$  is the largest,

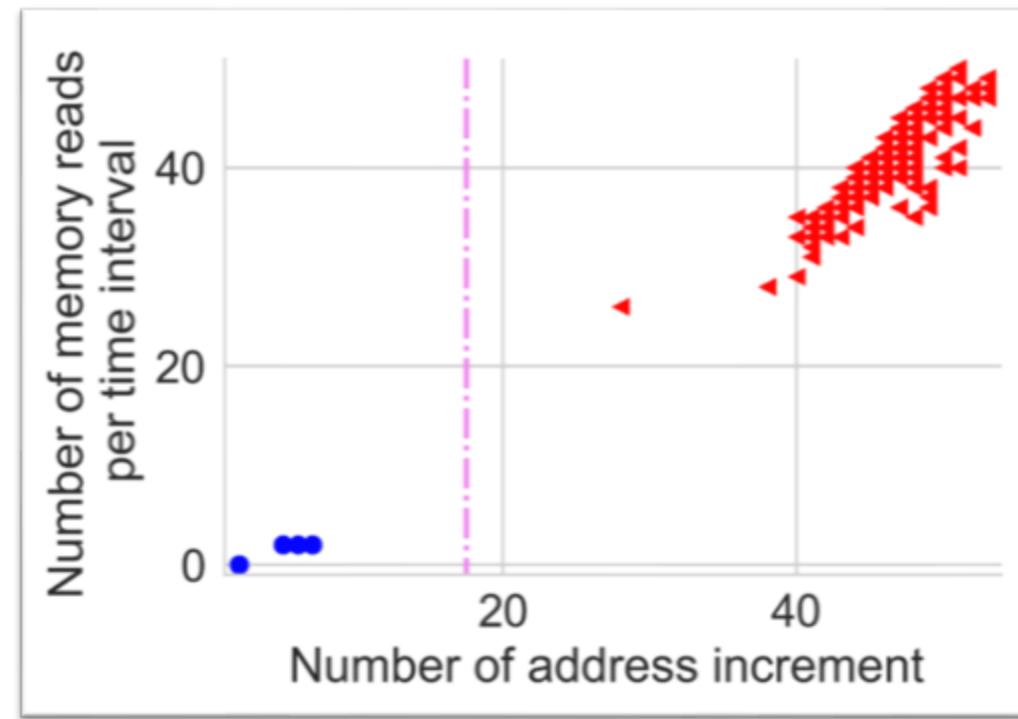
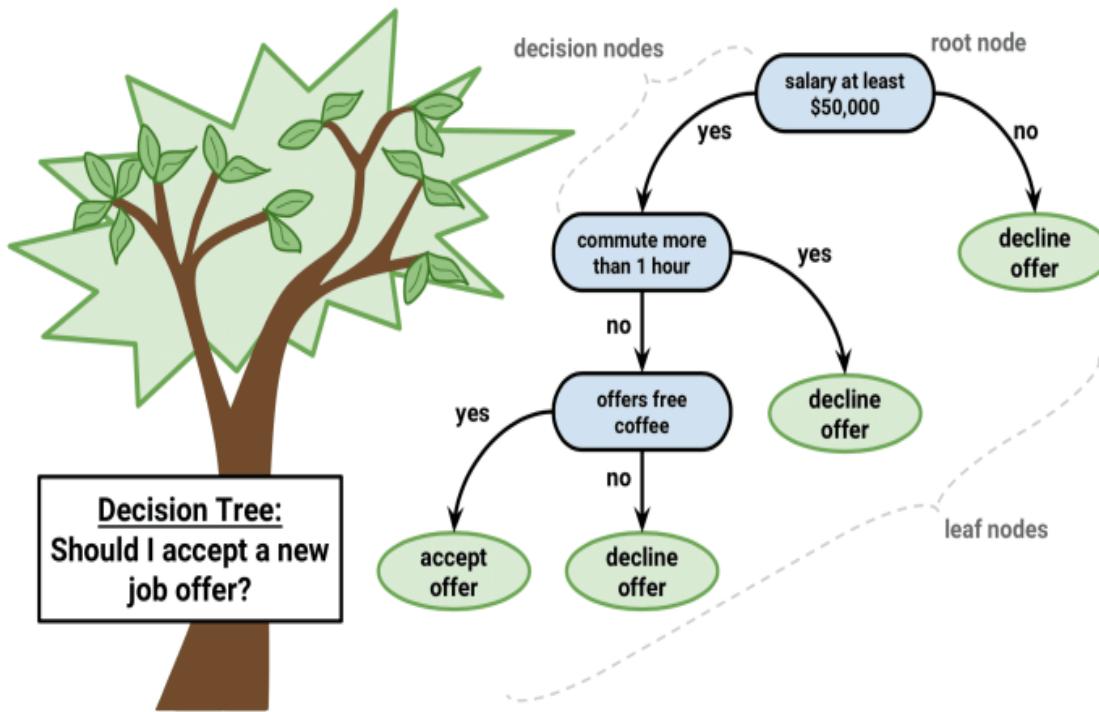


- **Assumption:**
  - Data have a notion of distance (data are in a metric space)
- **Intuition:**
  - Lazy learner → store all the training data and for every new incoming new observation, the algorithm will try to find the k nearest neighbor and do a majority voting



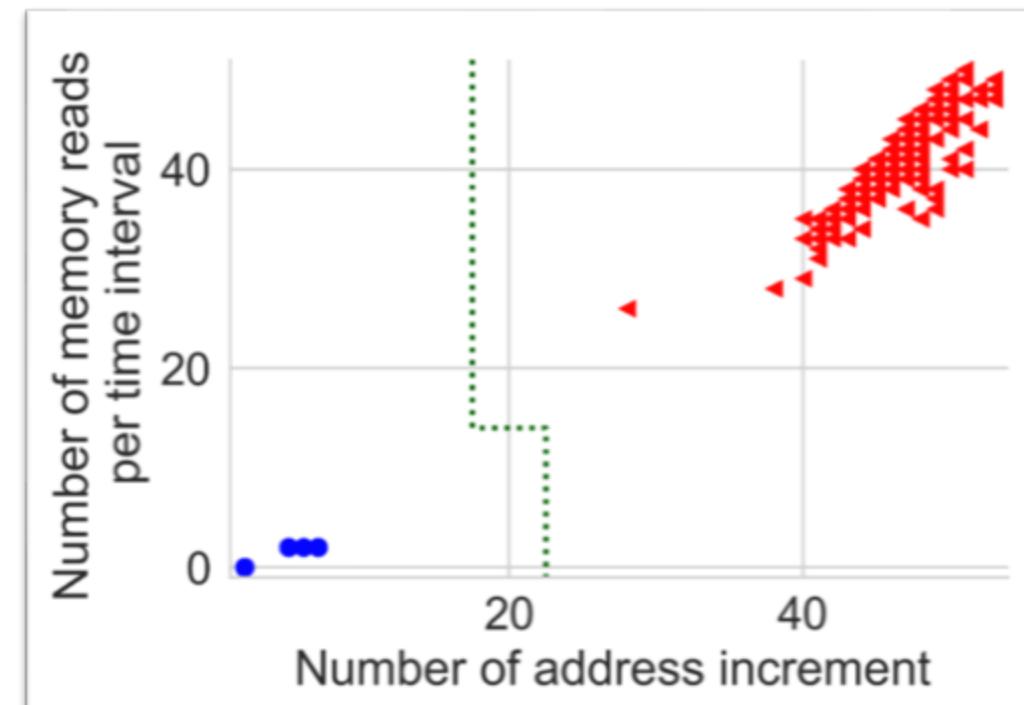
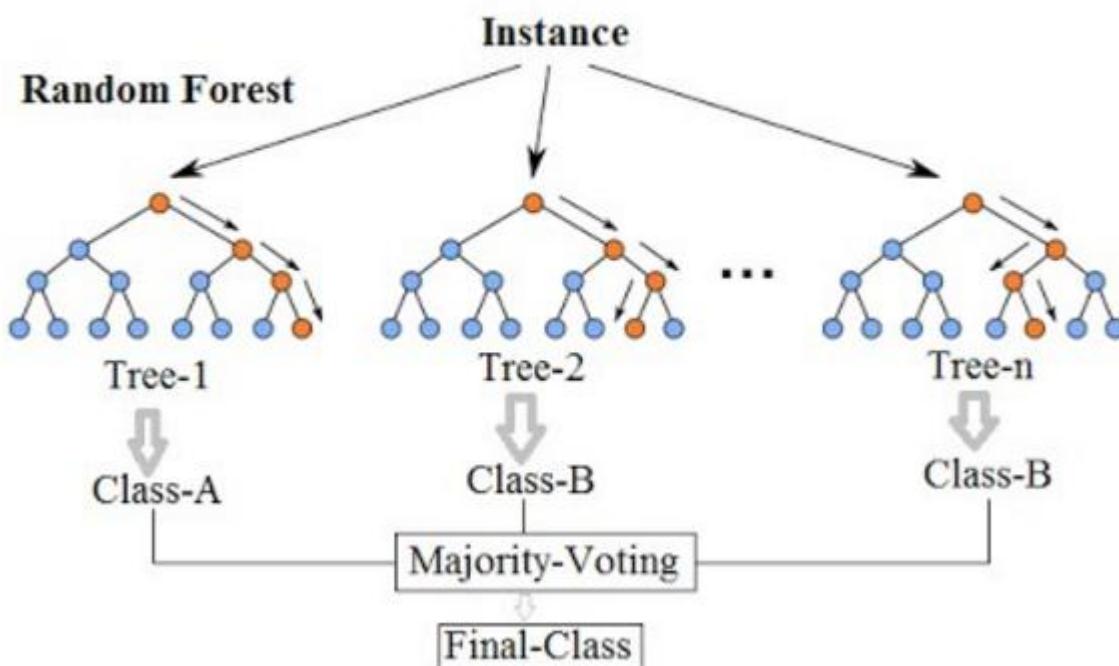
# DECISION TREE

- Assumption :
  - None
- Intuition
  - Decompose a complex decision into a union of several simpler decision

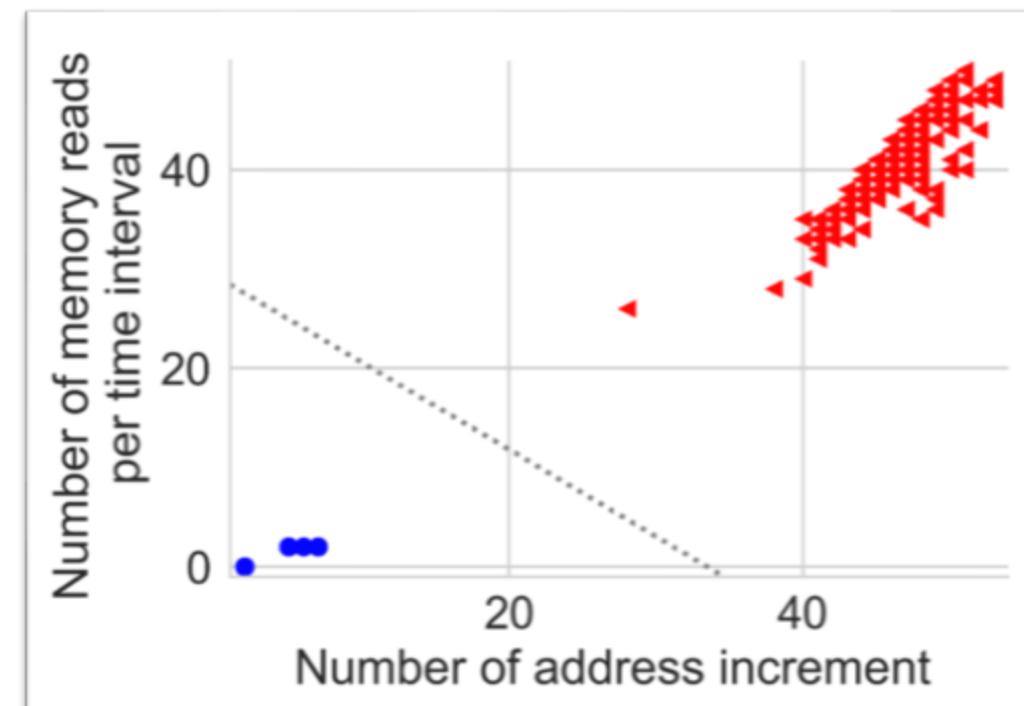


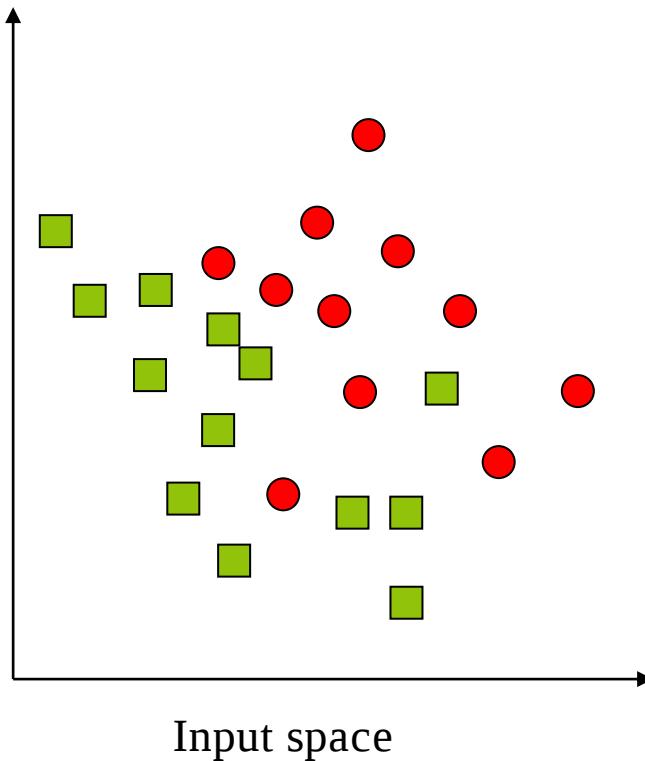
- Assumption :
  - None
- Intuition
  - a collection or ensemble of simple tree predictors, each capable of producing a response when presented with a set of predictor values.

### Random Forest Simplified

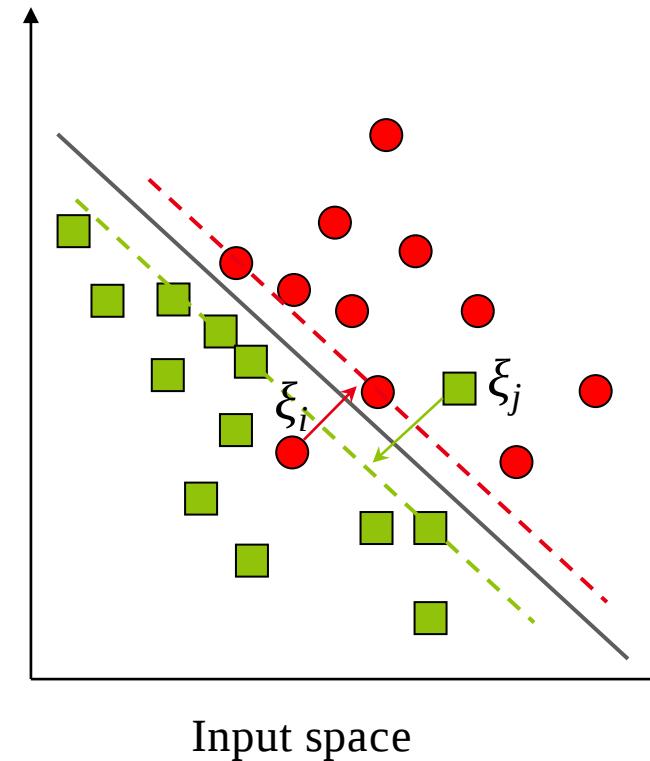


- **Assumption**
  - Linear SVM: the decision boundary is linear
- **Intuition:**
  - The decision boundary should be as far away from the data of both classes as possible
    - ➔ We should maximize the margin  $m = \frac{1}{\|w\|}$
  - This maximum-margin separator is determined by a subset of the data points (support vectors).
    - ➔ It will be useful computationally if only a small fraction of the data points are support vectors, because we use the support vectors to decide which side of the separator a test case is on.

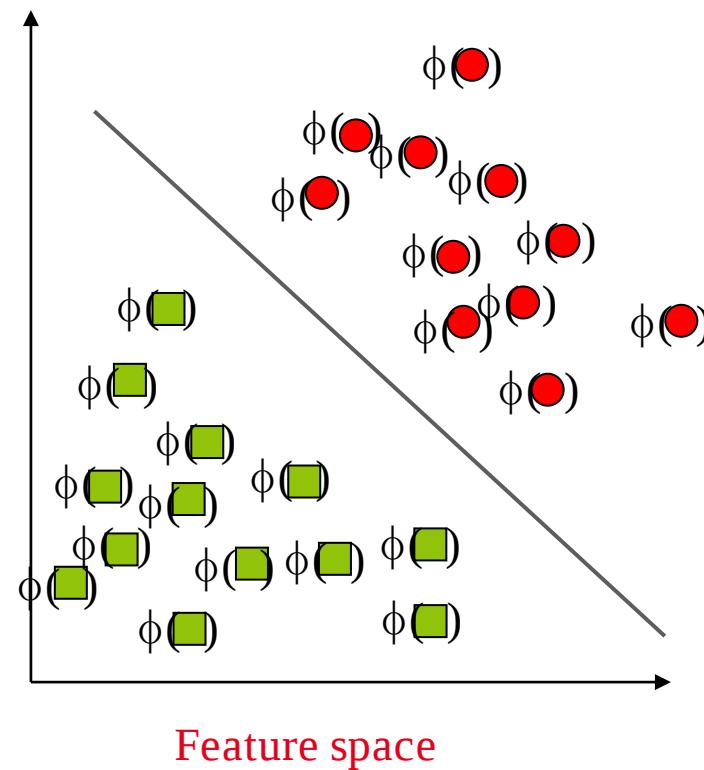
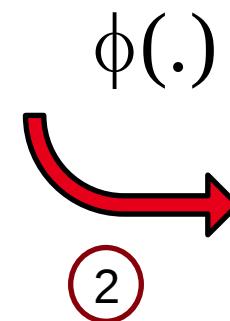
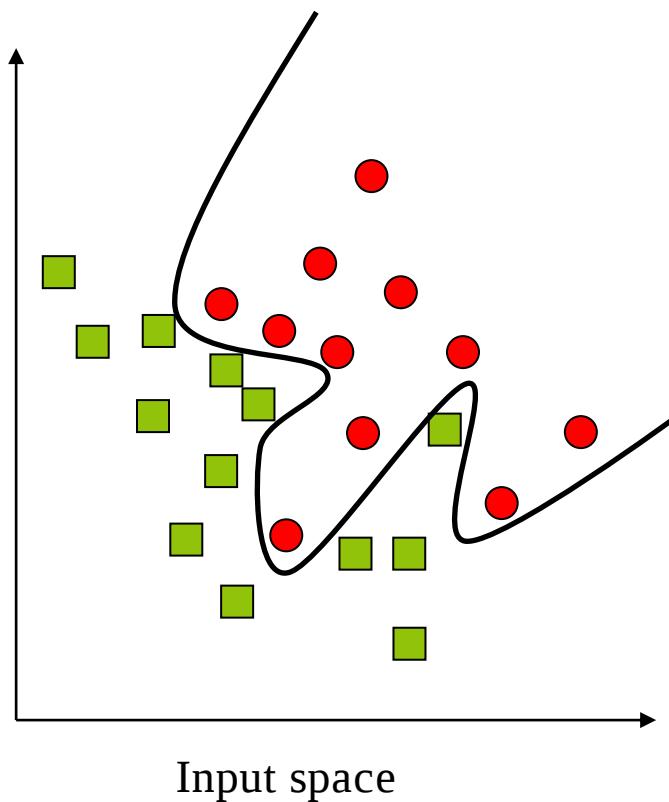




① Slack variables  $\xi_i$  can be added to allow misclassification of difficult or noisy examples.



Projection to higher dimensional space where we can find a linear separator

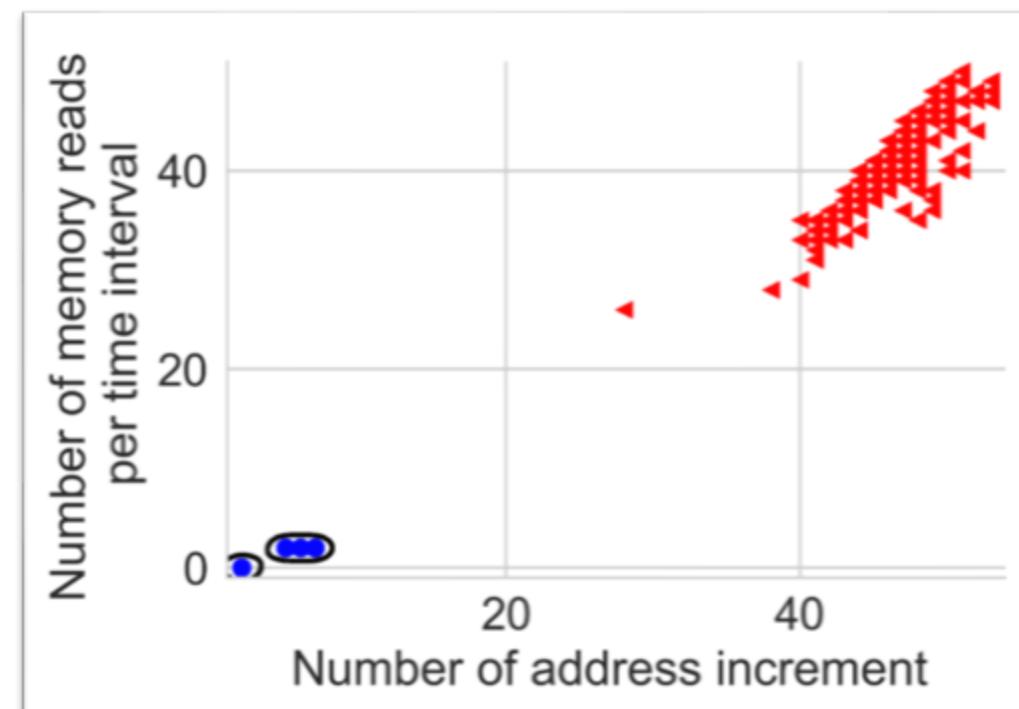


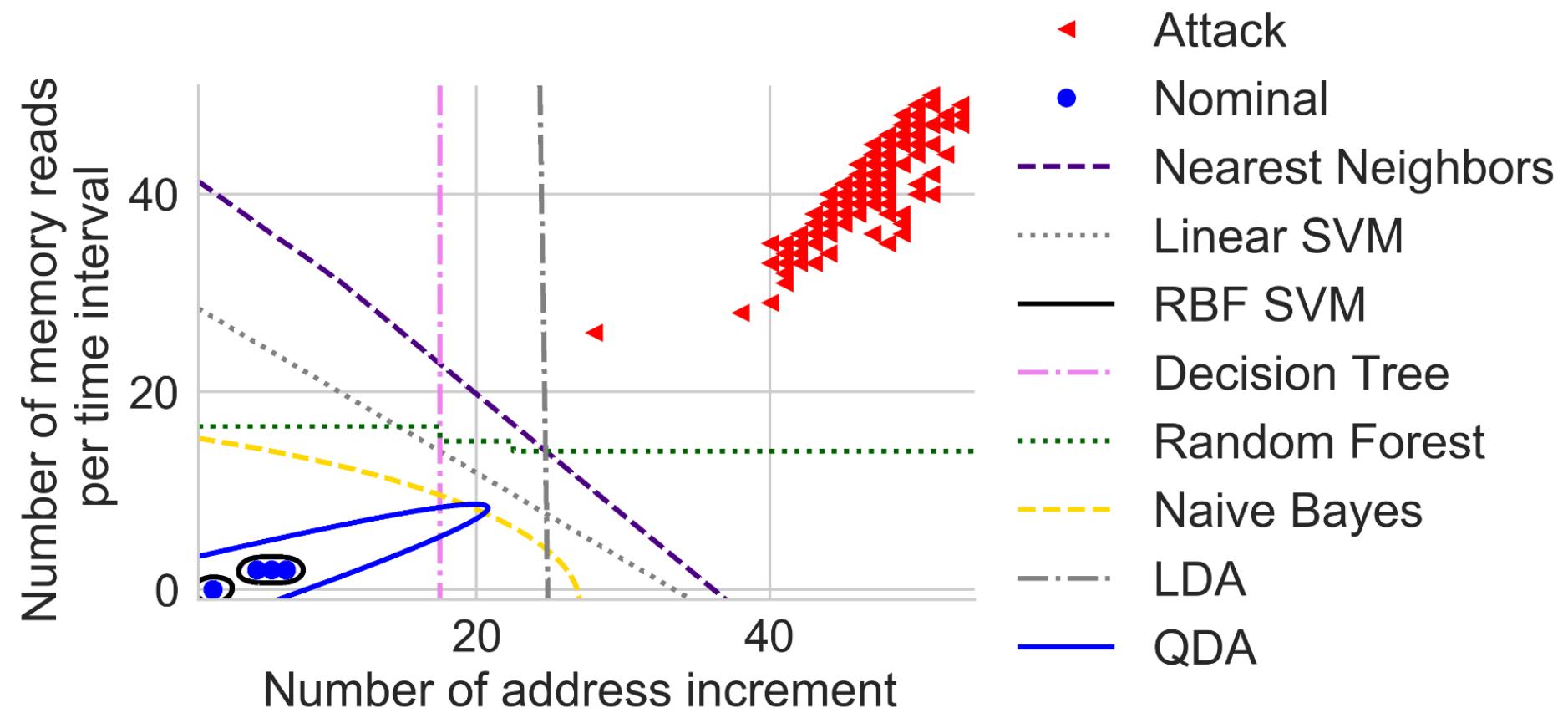
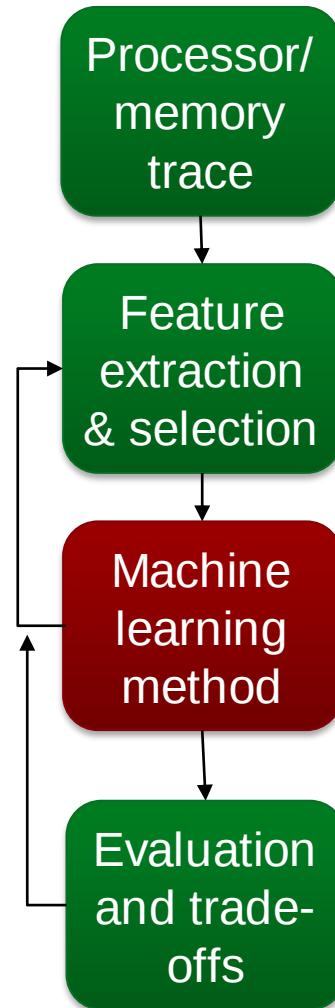
- The final classification rule is quite simple:

$$f(x_{tes}) = \text{sig}(\eta b + \sum_{s \in SV} \alpha_s y_s K(x_{tes}, x_s))$$

The set of support vectors      Lagrange parameter

- All the cleverness goes into selecting the support vectors that maximize the margin and computing the weight to use on each support vector.
- We also need to choose a good kernel function and set the parameters of the used kernel
- Popular kernels :
  - polynomial of a degree d  $K(x_i, x_j) = (x_i^T x_j + 1)^d$
  - radial basis function  $K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$





## EVALUATION METRICS

Processor/  
memory  
trace

Feature  
extraction  
& selection

Classifiers

Evaluation  
and trade-  
offs

- **False positive rate:** number of false alarms generated by the classifier

$$FP = \frac{R}{FP + TN}$$

- **False negative rate:** number of miss detection by the classifier

$$FNR = \frac{FN}{FN + TP}$$

- **Precision:** is a measure of a classifiers exactness

$$PP = \frac{TP}{TP + FP}$$

- **Leakage:** number of bytes leaked before the classifier detects

- **Cost**

- Memory footprint of the classifier (in bytes)
- Computation (number of basic arithmetic operation needed to classify one instance)

## Principle

- In order to compare the computation cost in predicting the label of one instance for each classifier, we decomposed the learnt decision function of each classifier to basic arithmetic operations (additions, subtraction, comparison, multiplications, square root; exponential and divisions)
- The memory cost is computed by calculating the number of variables needed by each classifier

Classifier	Add(+)/Sub(-) /comparison	Mul (x)	Sqrt()	Exp()	Div
LSVM	$(d + 1)n_s - 1$	$(d + 2)n_s$	0	0	0
RSVM	$(d + 1)n_s$	$(d + 3)n_s$	0	$n_s$	$n_s$
KNN	$2n(d + 1) - 2 \times k$	$n \times d$	$n$	0	0
LDA	$d$	$d$	0	0	0
QDA	$d^2 + d$	$d^2 + 2 d$	0	0	0
Naïve Bayes	$n_c$	2d	$d$	$d$	$2d$
Random Forest	$n_{tre} (h^e + 1)$	0	0	0	0
Decision Tree	$h$	0	0	0	0

d: number of features

n: number of observations in training dataset

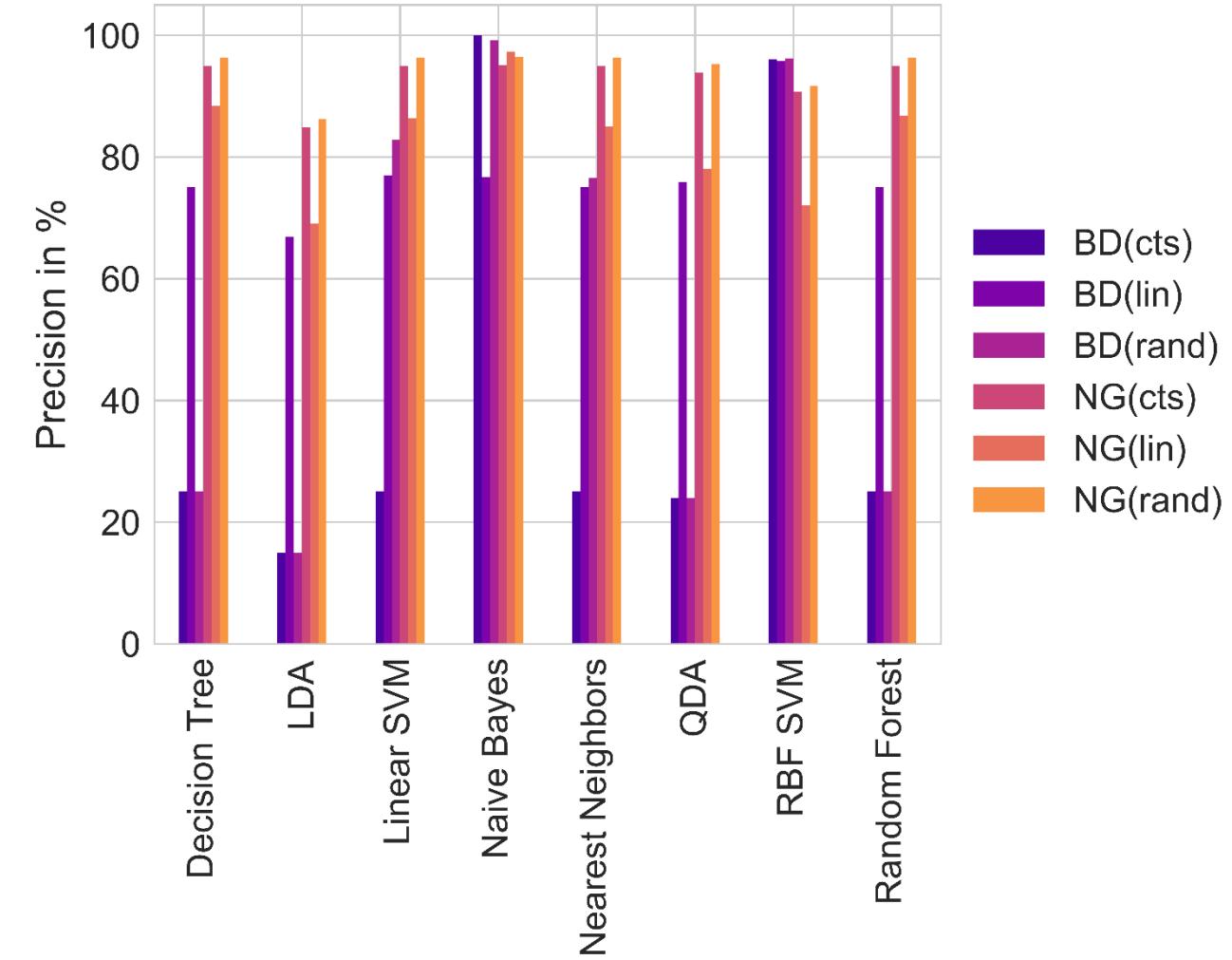
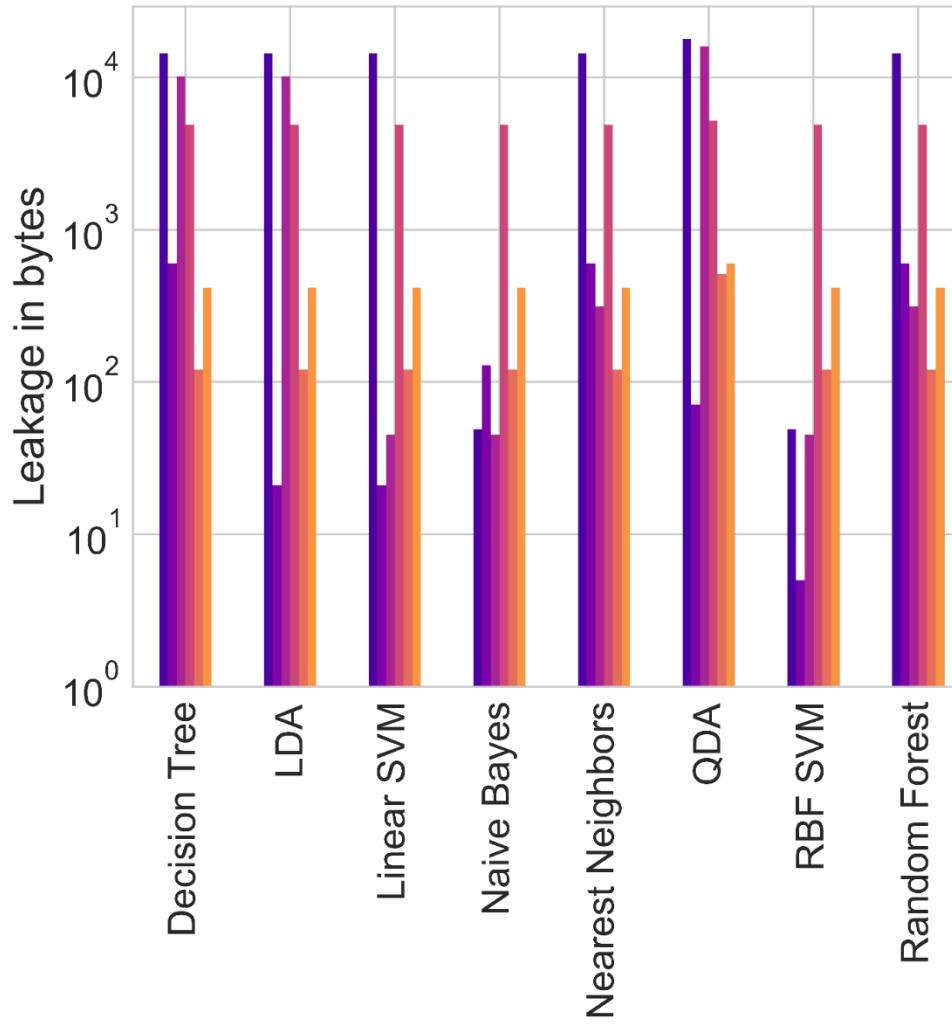
$n_c$ : number of classes

$n_{tree}$ : number of trees in the random forest

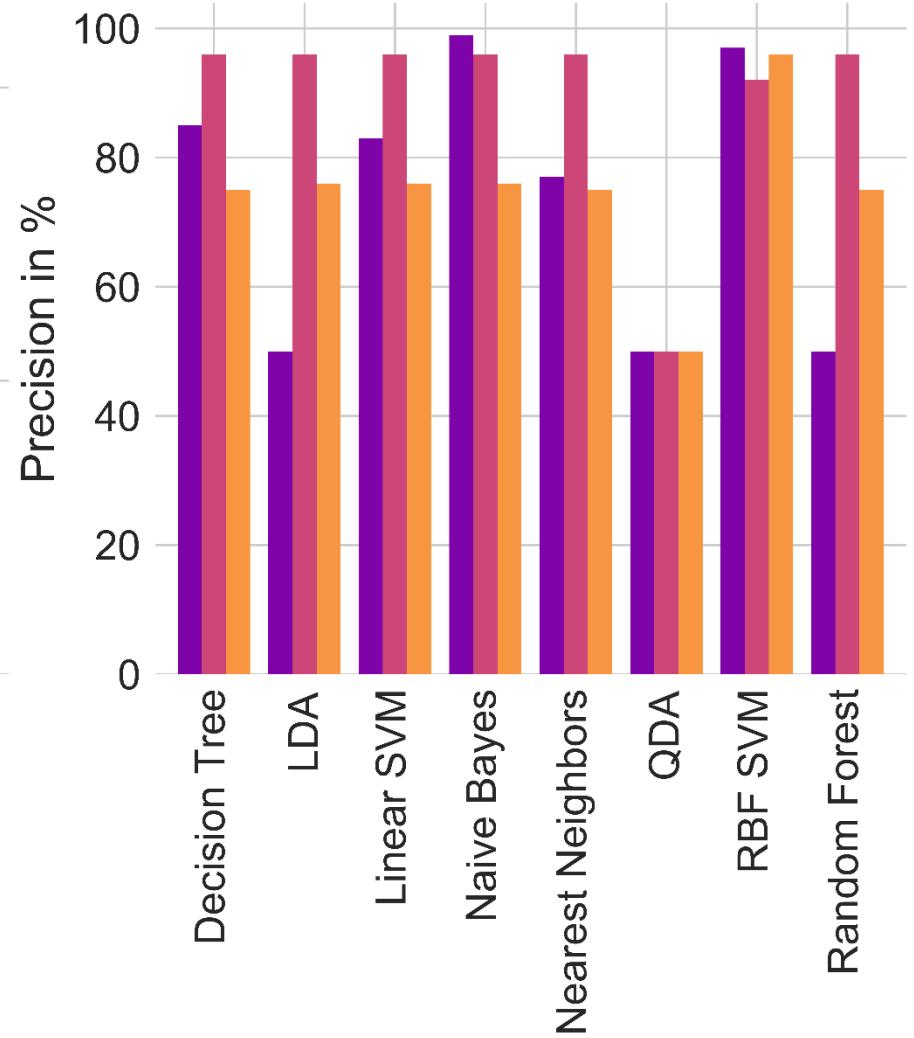
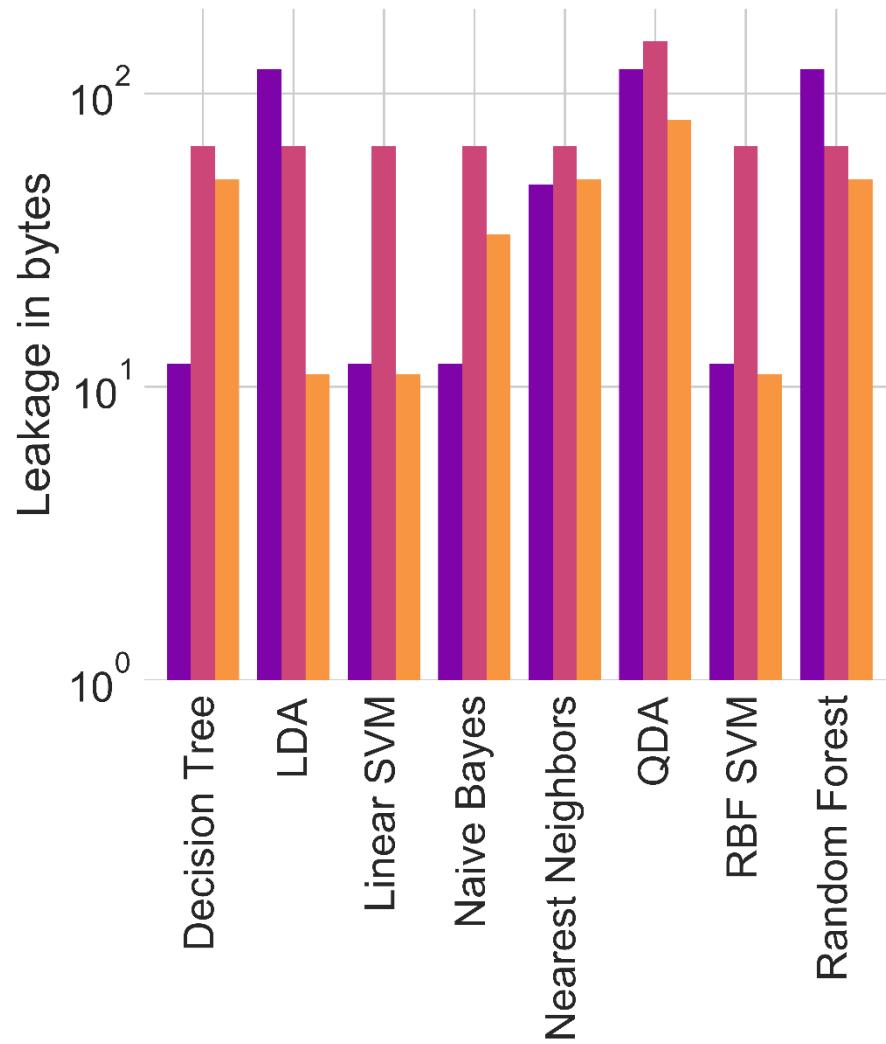
$n_s$ : number of support vectors

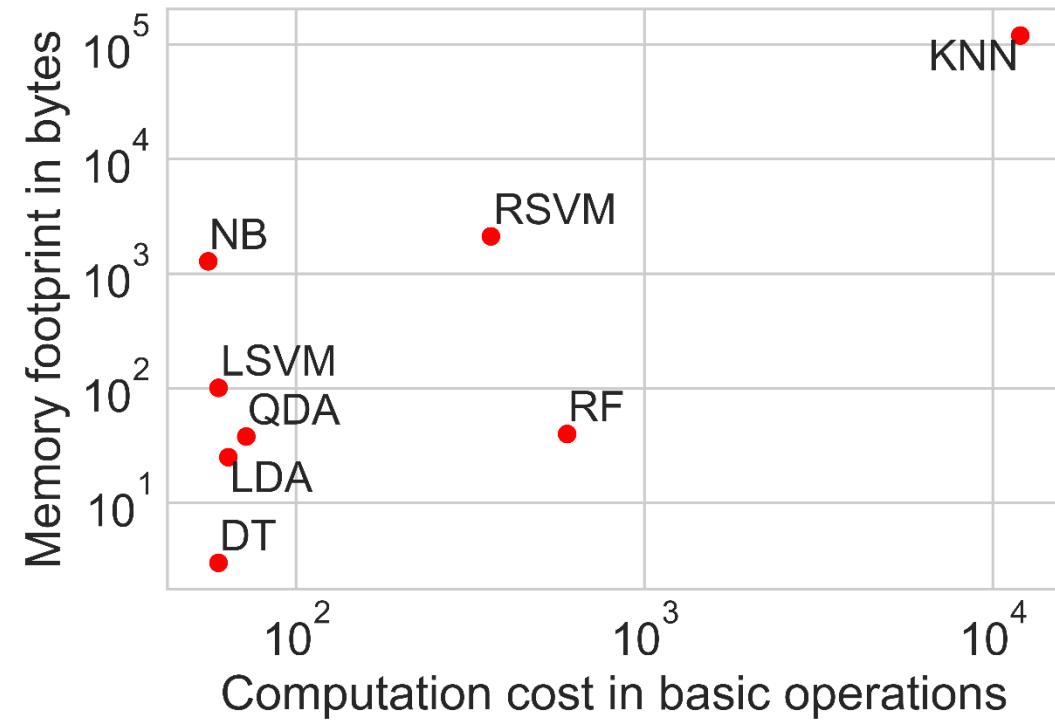
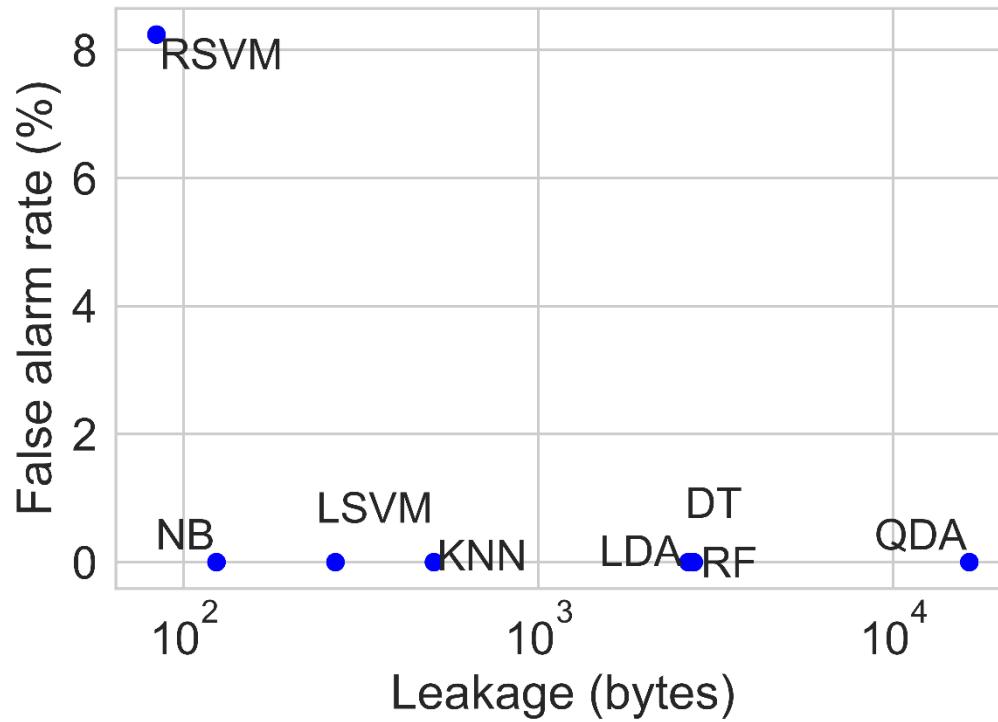
k: number of neighbors

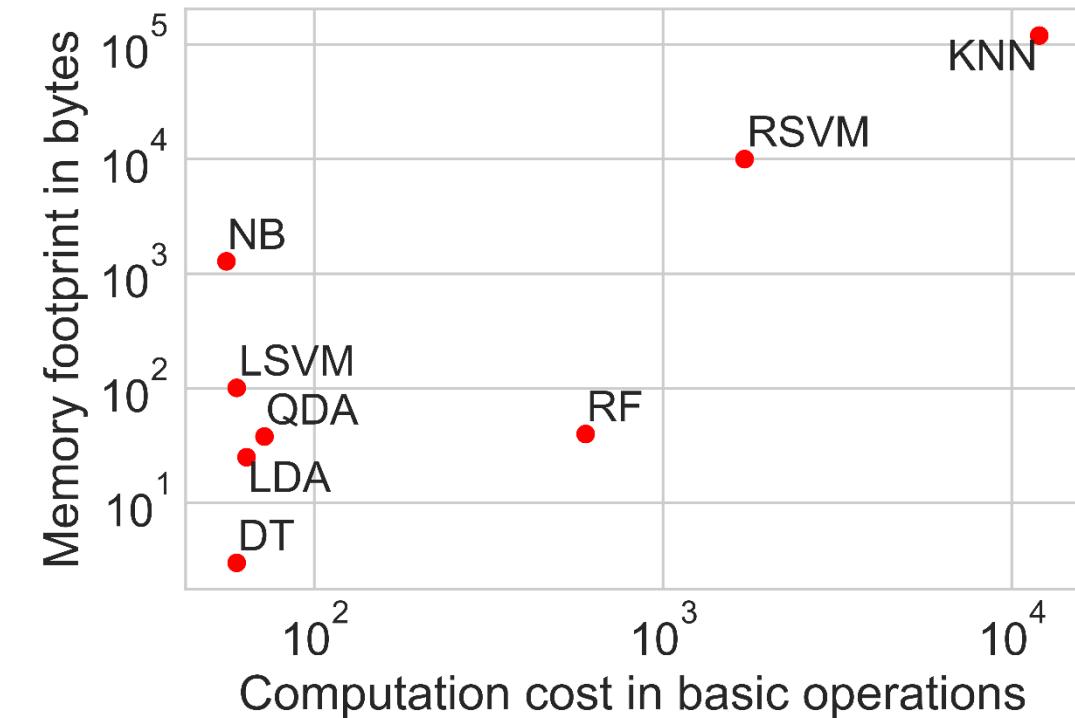
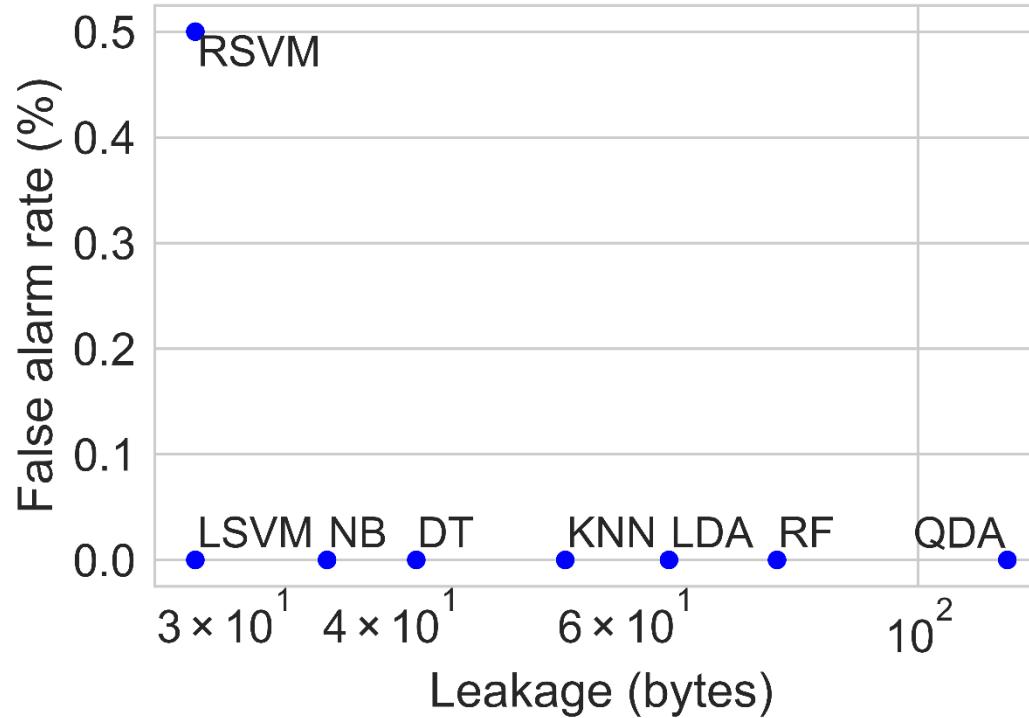
h: depth of the tree

DETECTION PRECISION & LEAKAGE OF CLASSIFIERS TRAINED  
ON CLASSIC DUMP

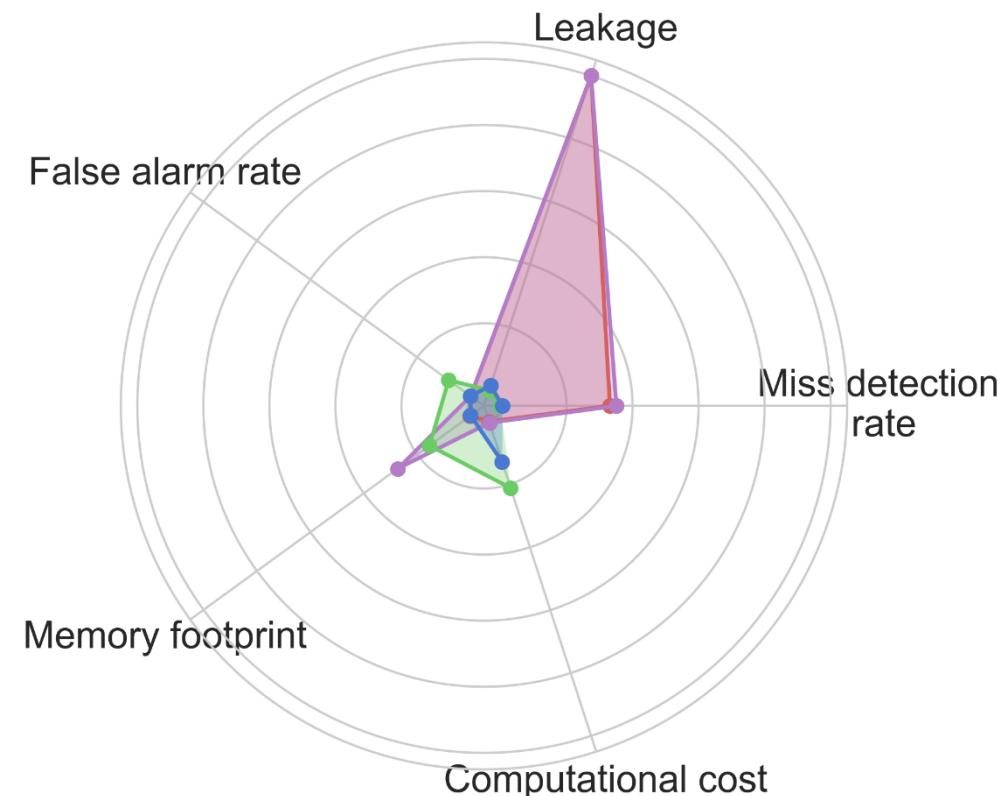
## DETECTION PRECISION &amp; LEAKAGE OF CLASSIFIERS TRAINED ON VARIANTS DUMPS



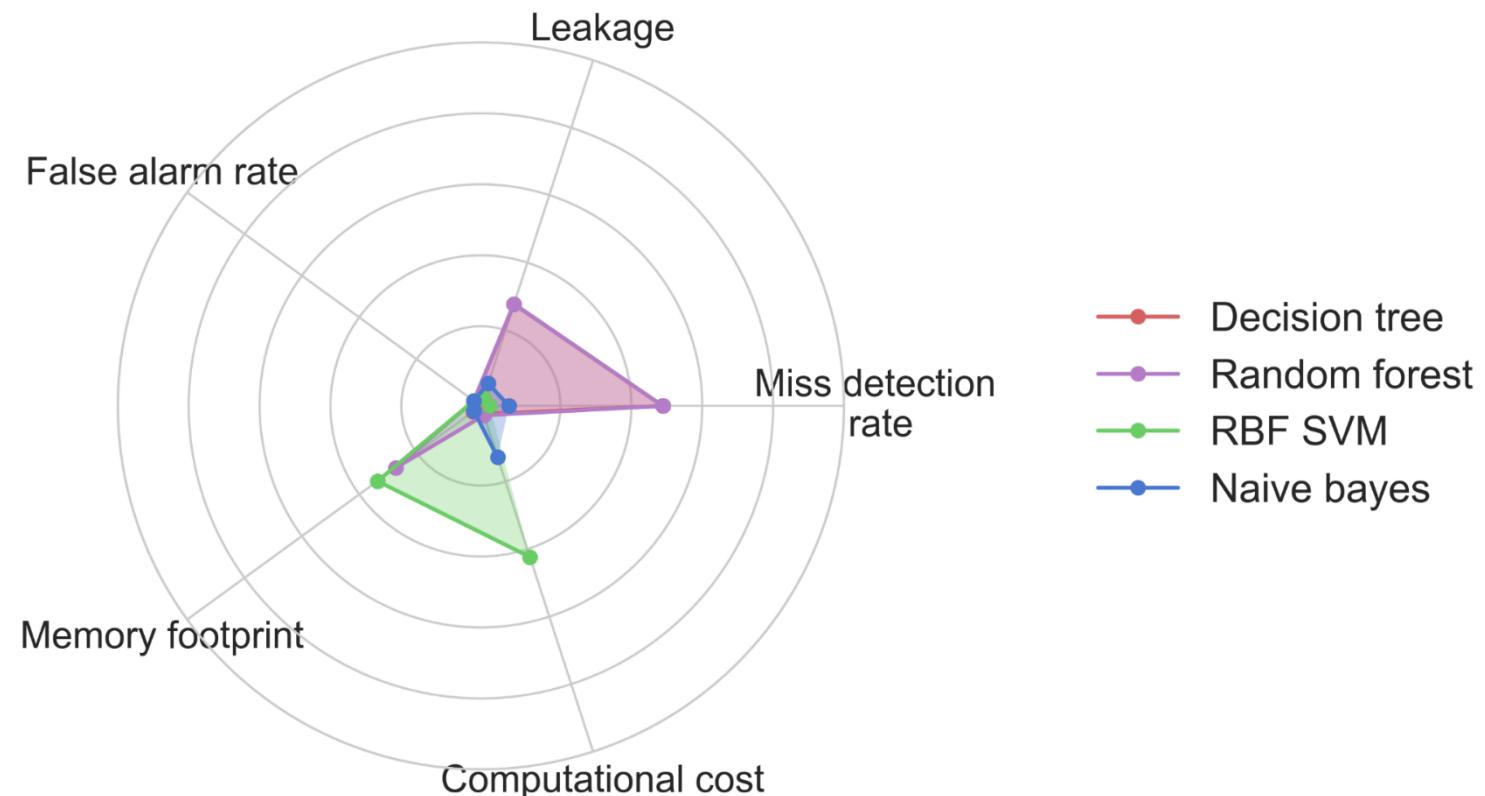




## COMPARISON OF CLASSIFIERS PERFORMANCE



Trained on classic dump



Trained on dump variants

- **Take out:**
  - Binary classifiers are a great choice for low cost detectors
  - Diversifying the training dataset can increase the accuracy of the detection but that comes at the cost of the implementation complexity
  - Even when trained on limited examples of attacks binary classifiers were able to detect efficiently (few bytes leakage and detection accuracy around 90%)
- **Next step**
  - Implementation on hardware
  - Exploration of other types of attacks
  - Evaluation of mimicry attack cost to evade the detection

# Q&A

Contact us for more details:  
[Sanaa.kerroumi@cea.fr](mailto:Sanaa.kerroumi@cea.fr)  
[Anca.Molnos@cea.fr](mailto:Anca.Molnos@cea.fr)  
[Damien.Couroussé@cea.fr](mailto:Damien.Couroussé@cea.fr)

