

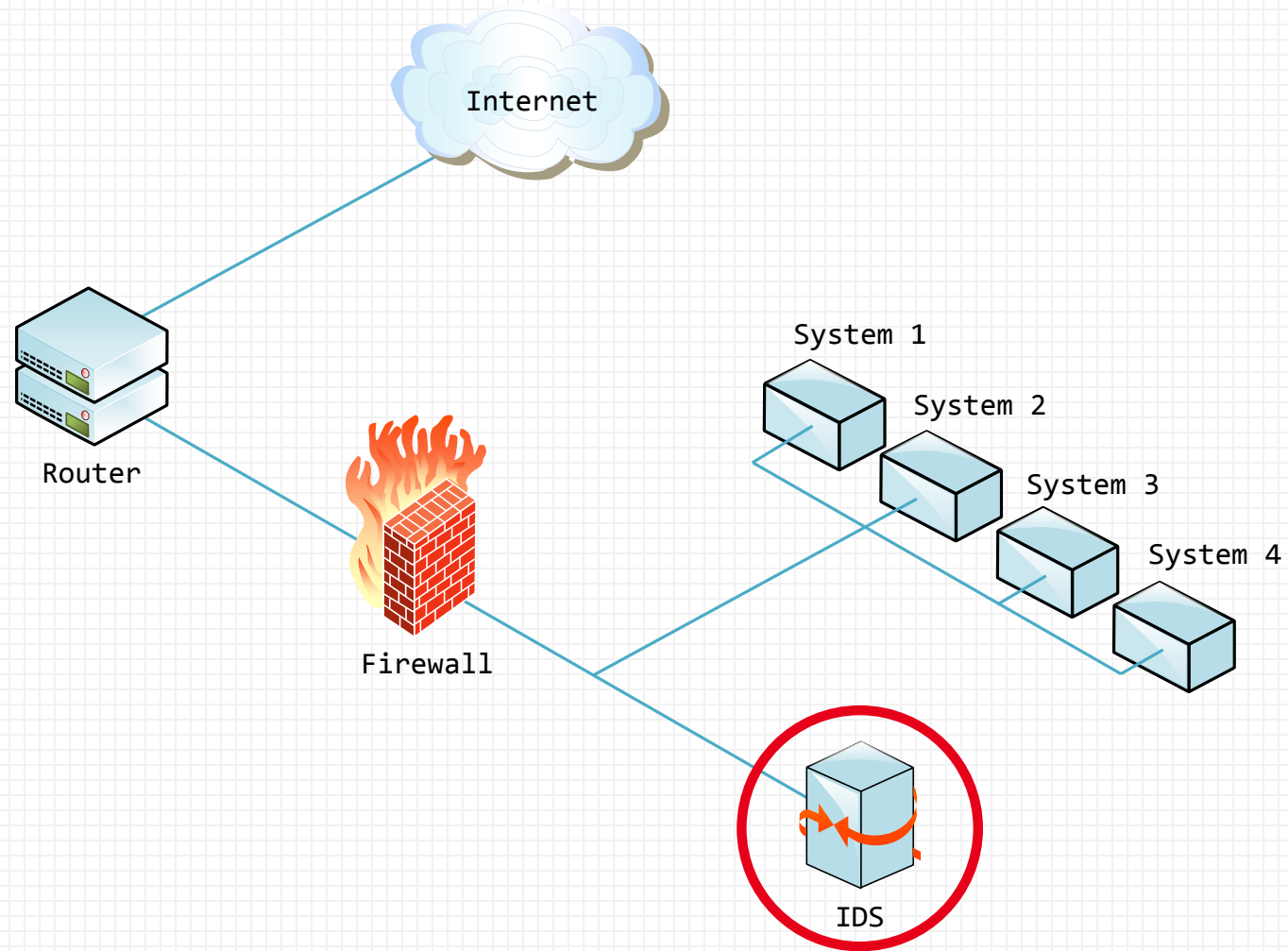


AUTOMATISATION DU PROCESSUS D'ENTRAÎNEMENT POUR LA DÉTECTION D'INTRUSION



- **Intrusion Detection Systems**
- **Problem refinement**
 - NSL-KDD dataset
 - Classification process
- **Preprocessing**
- **Classifiers**
- **Ensemble learning**
- **Results**

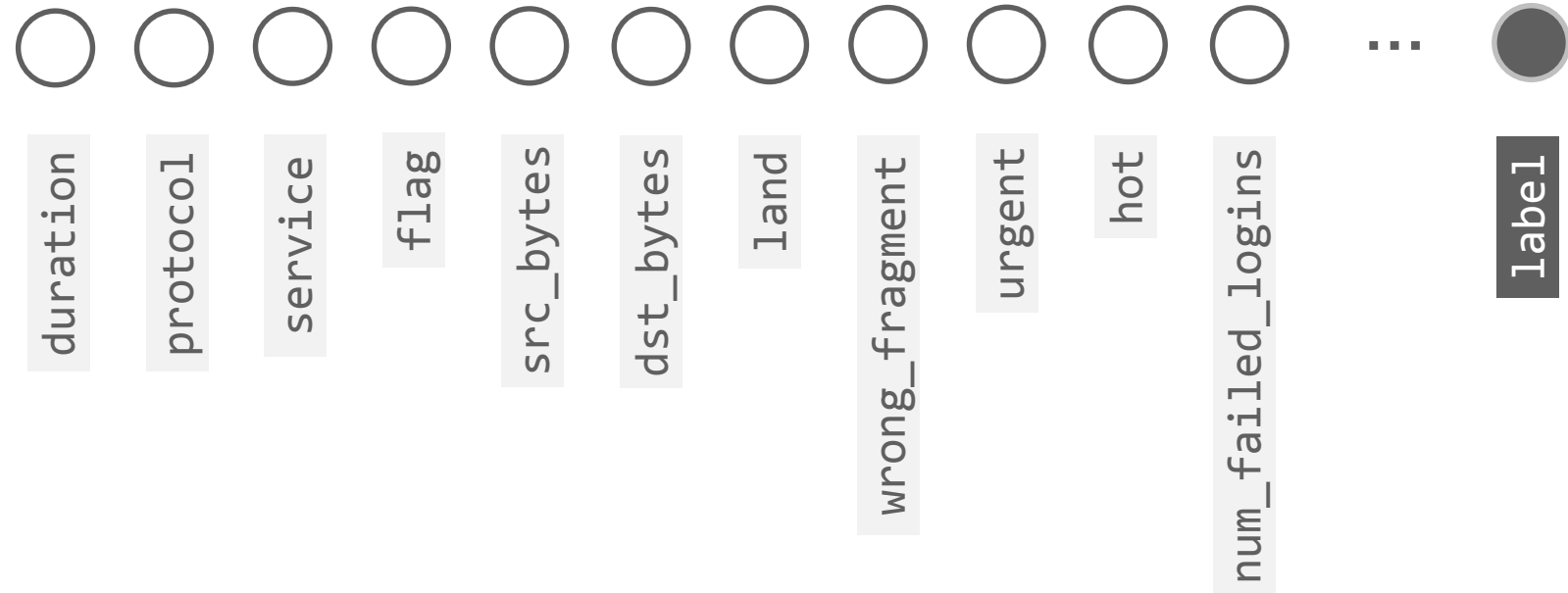
INTRUSION DETECTION SYSTEMS



NSL-KDD DATASET



NSL-KDD DATASET



41 features + 1 label

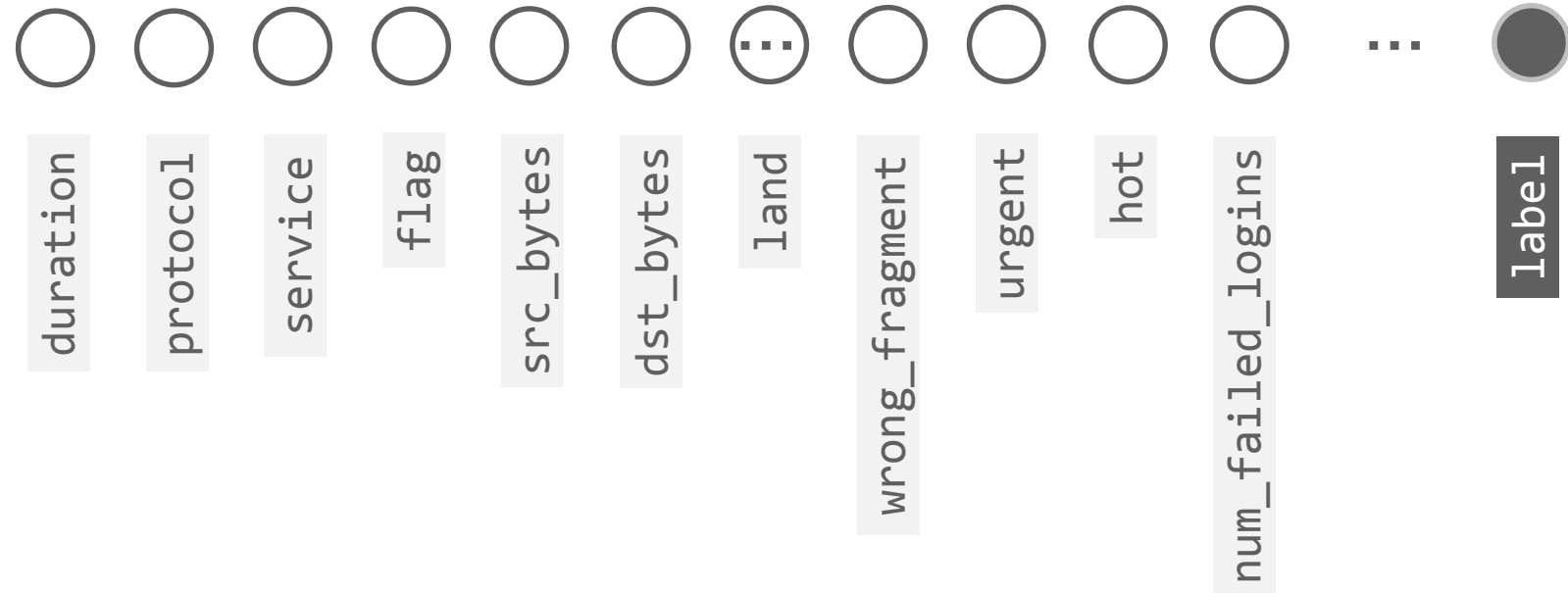
Example 1:

2	tcp	ftp	SF	491	0	0	0	0	0	0	0	...	normal
---	-----	-----	----	-----	---	---	---	---	---	---	---	-----	--------

Example 2:

0	tcp	priv	REJ	0	0	0	0	0	0	0	0	...	neptune
---	-----	------	-----	---	---	---	---	---	---	---	---	-----	---------

CLASSIFICATION PROCESS



41 features + 1 label

Example 1:

2	tcp	ftp	SF	491	0	0	0	0	0	0	...	normal
---	-----	-----	----	-----	---	---	---	---	---	---	-----	--------

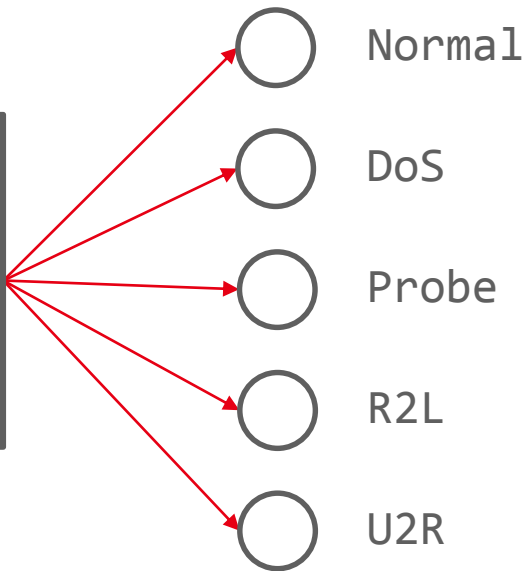
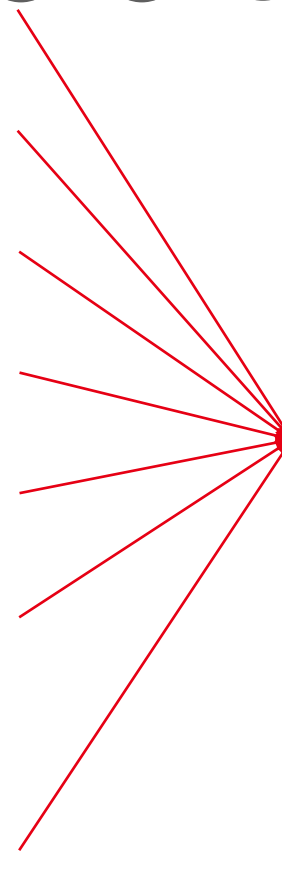
Example 2:

0	tcp	priv	REJ	0	0	0	0	0	0	0	...	neptune
---	-----	------	-----	---	---	---	---	---	---	---	-----	---------

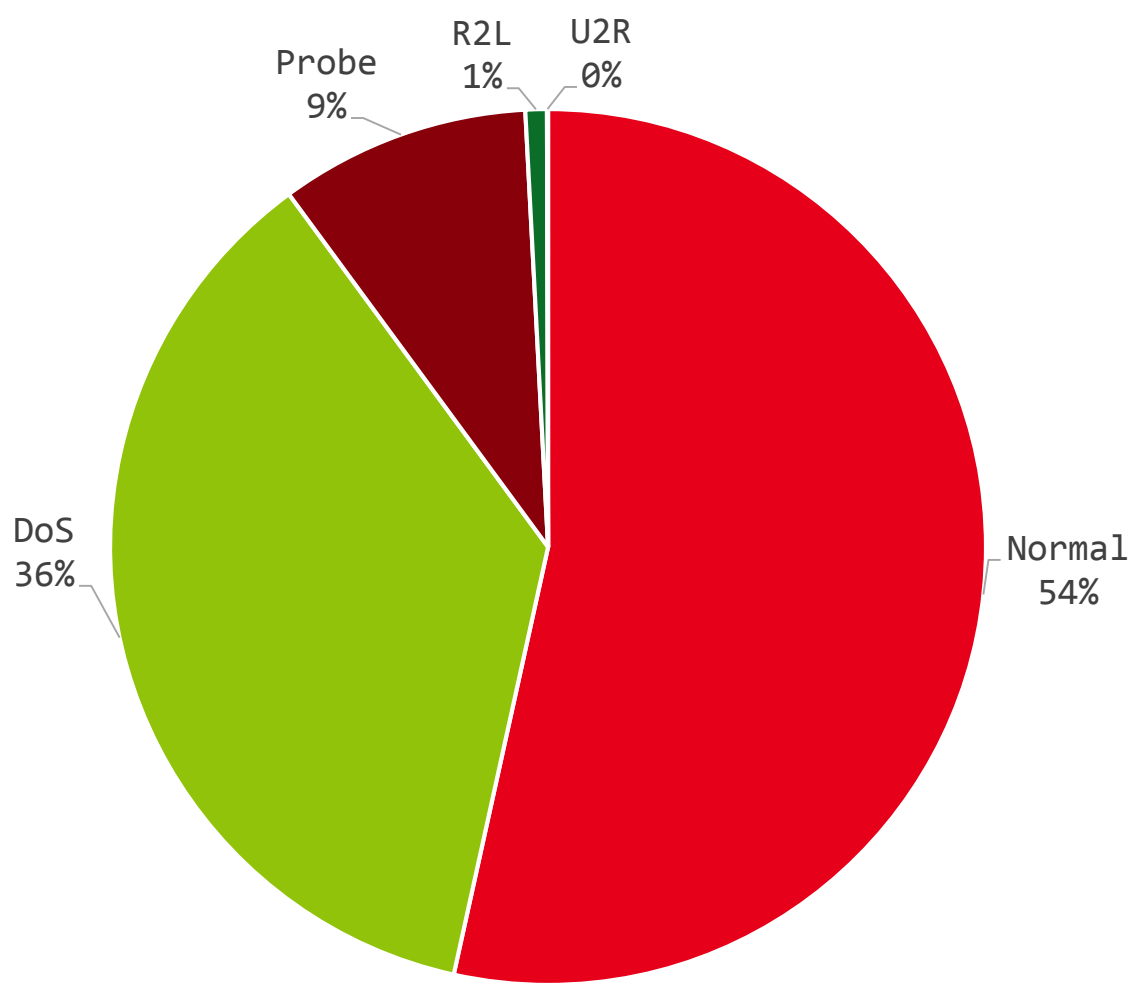
CLASSIFICATION PROCESS

NSL-KDD

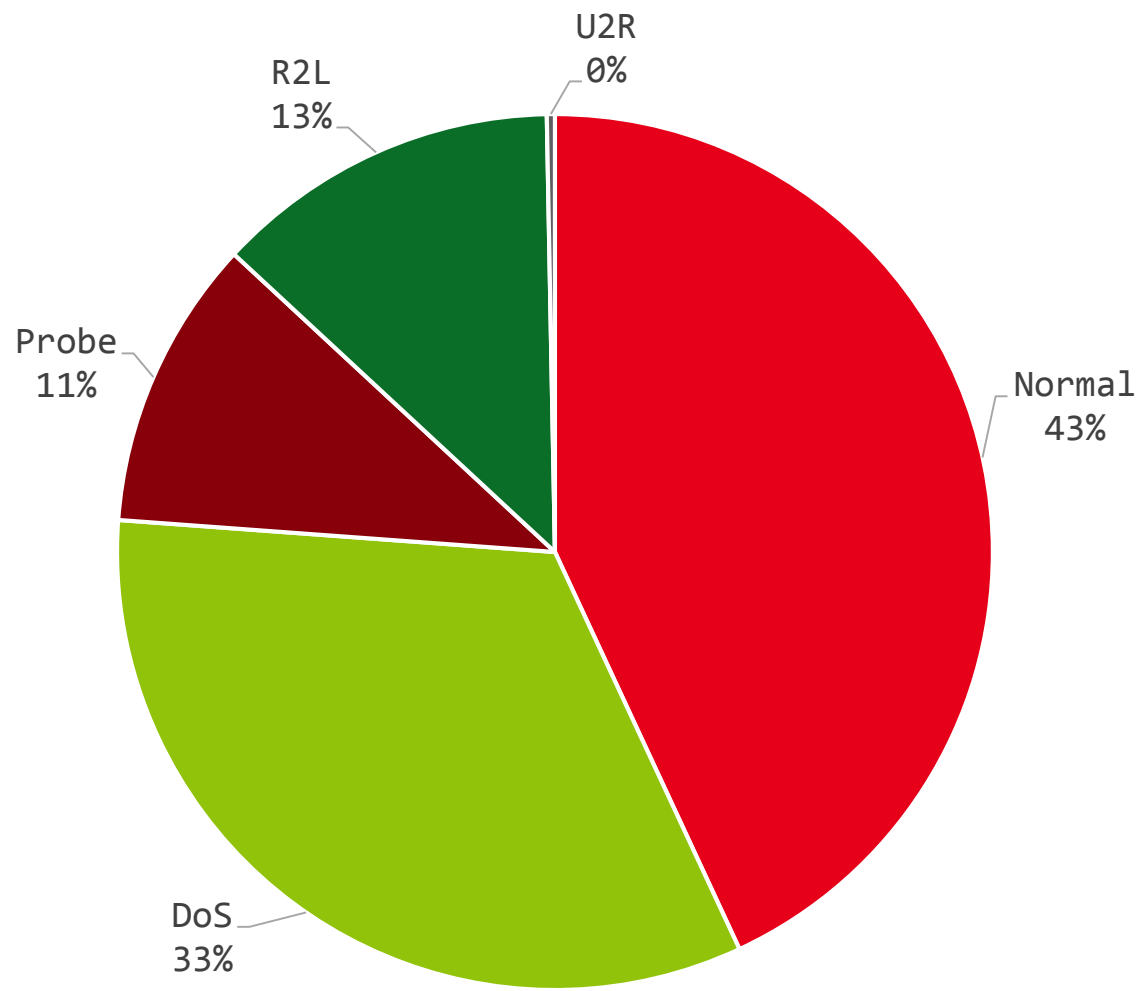
41 features



PREPROCESSING

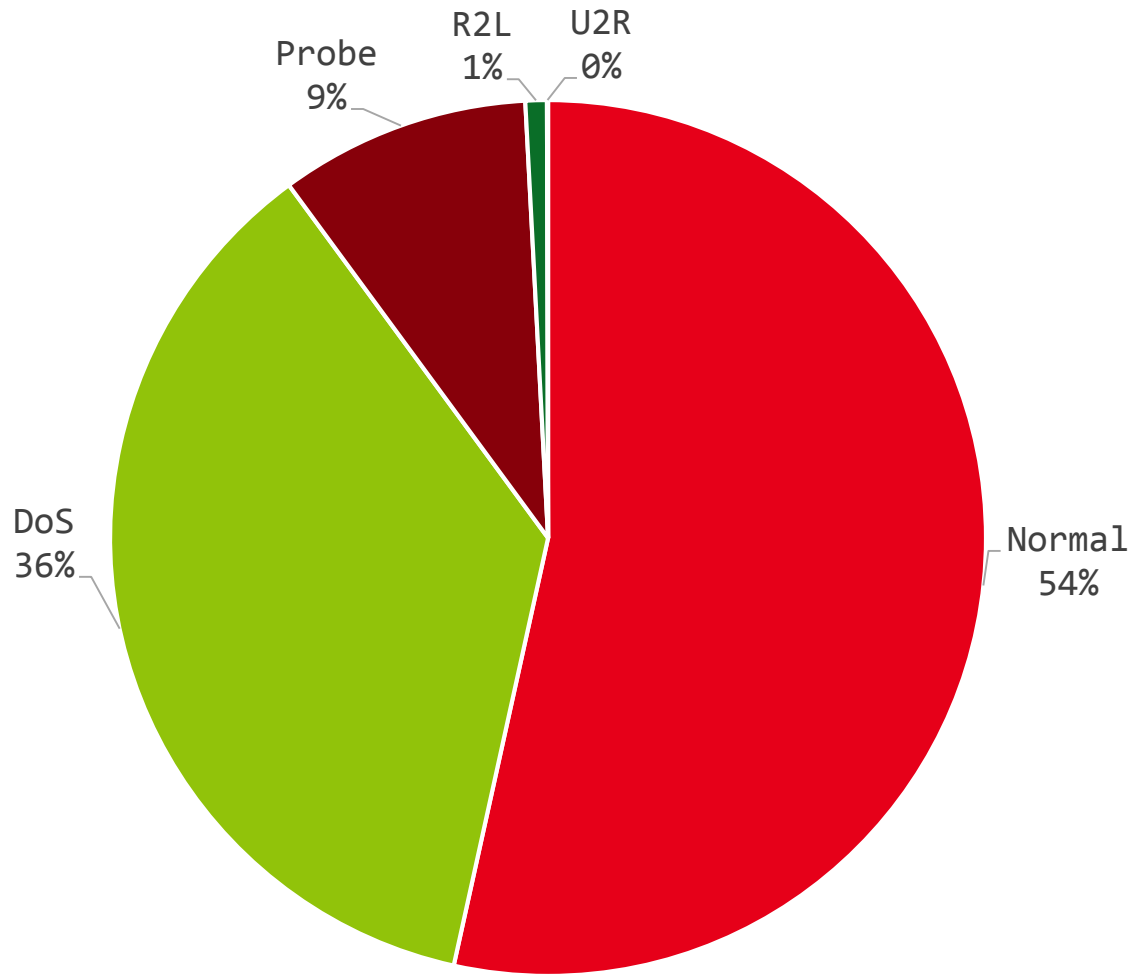


Training set



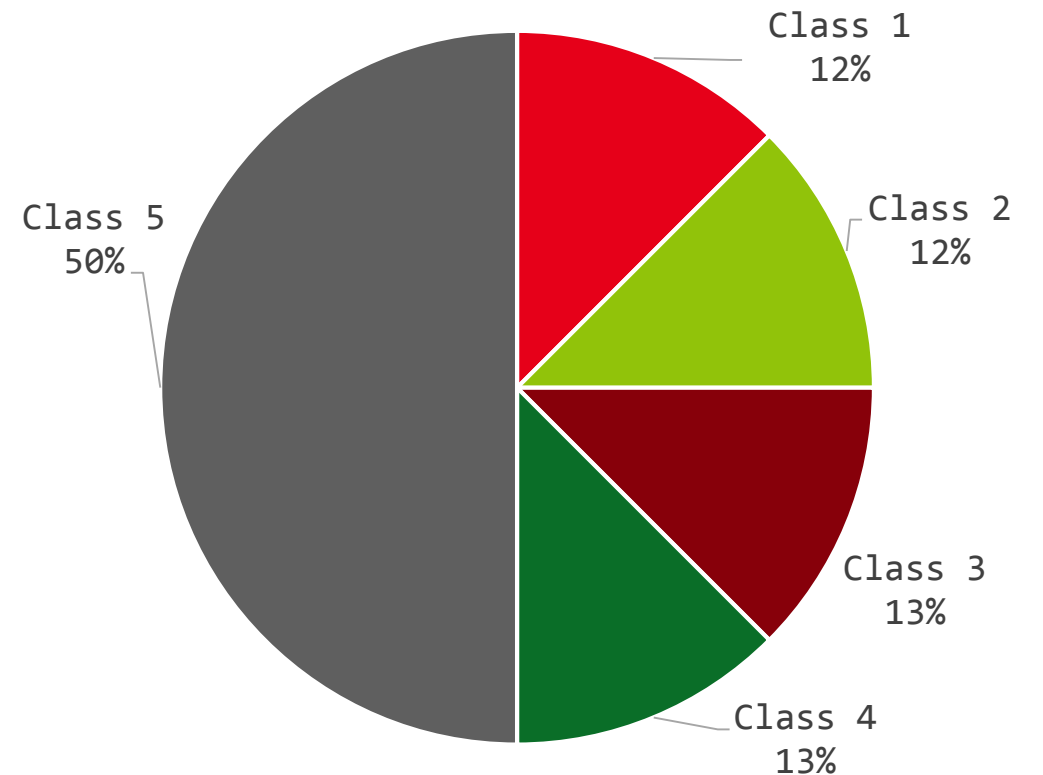
Test set

PREPROCESSING



Training set

Each class is successively overrepresented to create 5 specialized training sets



Specialized training + validation set

ALGORITHMS

- *No sampling*
- *Cluster Centroids*
- *Random Under Sampler*
- *NearMiss*
- *Edited Nearest Neighbours*
- *Repeated Edited Nearest Neighbours*
- *Condensed Nearest Neighbour*
- *One Sided Selection*
- *AIKNN*
- *Instance Hardness Threshold*
- *Random Over Sampler*
- *SMOTE*
- *Borderline 1 SMOTE*
- *Borderline 2 SMOTE*
- *SVM SMOTE*
- *ADASYN*
- *SMOTEENN*
- *SMOTE Tomek*

Under-sampling

Over-sampling

We want to **combine** the best **under-sampling** algorithm with the best **over-sampling** algorithm

Best algorithm = best ROC AUC score

(classification with a simple MLP)

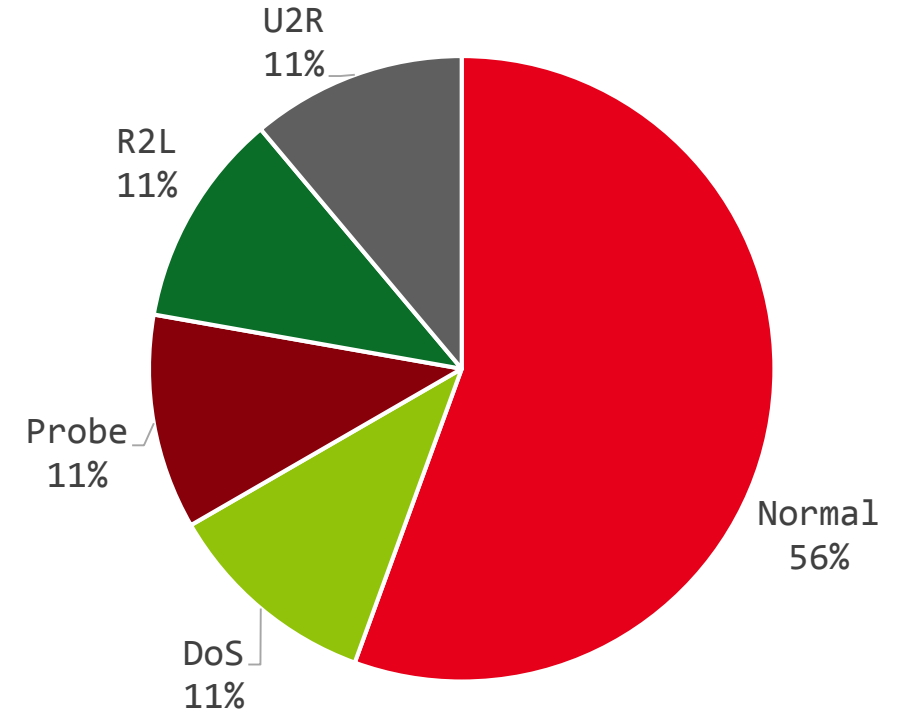
Under-sampling

Over-sampling

ALGORITHMS

ROC AUC SCORE

	Normal
• <i>No sampling</i>	0.9488
• <i>Cluster Centroids</i>	0.9378
• <i>Random Under Sampler</i>	0.9616
• <i>NearMiss</i>	0.9377
• <i>Edited Nearest Neighbours</i>	0.9382
• <i>Repeated Edited Nearest Neighbours</i>	0.9614
• <i>Condensed Nearest Neighbour</i>	0.9022
• <i>One Sided Selection</i>	0.9415
• AIKNN	0.9659
• <i>Instance Hardness Threshold</i>	0.6944
• <i>Random Over Sampler</i>	0.9611
• <i>SMOTE</i>	0.9423
• <i>Borderline 1 SMOTE</i>	0.9597
• <i>Borderline 2 SMOTE</i>	0.9542
• SVM SMOTE	0.9644
• <i>ADASYN</i>	0.9485
• <i>SMOTEENN</i>	0.9500
• <i>SMOTE Tomek</i>	0.9627



AIKNN + SVM SMOTE are used to:

- undersample *Normal*, *DoS* and *Probe*
- oversample *R2L* and *U2R*

U2R training + validation set

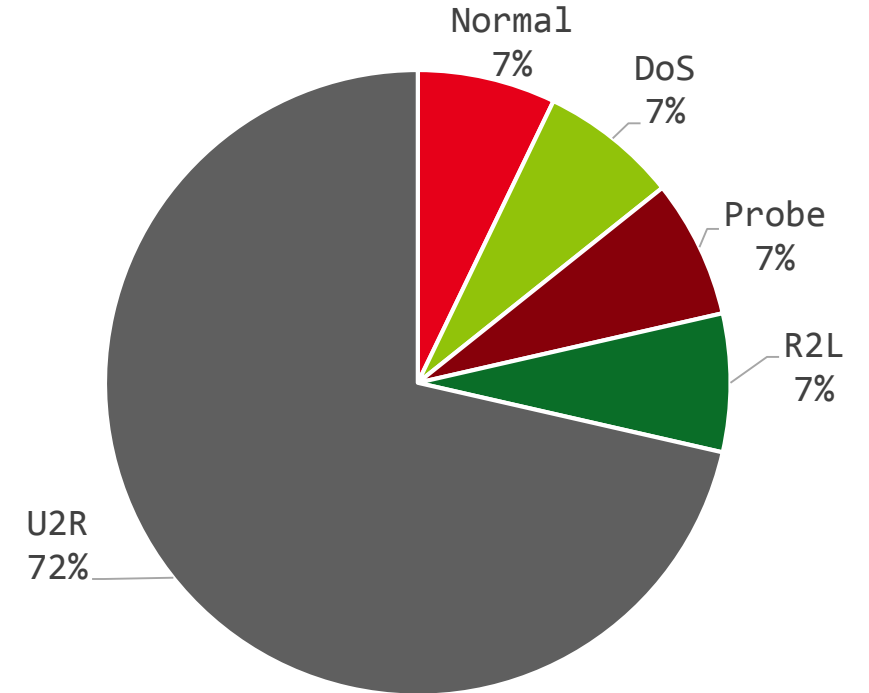
ALGORITHMS ROC AUC SCORE

Under-sampling

ALGORITHMS	ROC AUC SCORE
• <i>No sampling</i>	0.9605
• Cluster Centroids	0.9739
• <i>Random Under Sampler</i>	0.9715
• <i>NearMiss</i>	0.2762
• <i>Edited Nearest Neighbours</i>	0.9575
• <i>Repeated Edited Nearest Neighbours</i>	0.9626
• <i>Condensed Nearest Neighbour</i>	0.9174
• <i>One Sided Selection</i>	0.9435
• <i>AllKNN</i>	0.9701
• <i>Instance Hardness Threshold</i>	0.9458

Over-sampling

• <i>Random Over Sampler</i>	0.9714
• <i>SMOTE</i>	0.9511
• <i>Borderline 1 SMOTE</i>	0.9675
• <i>Borderline 2 SMOTE</i>	0.9601
• <i>SVM SMOTE</i>	0.9696
• <i>ADASYN</i>	0.9717
• <i>SMOTEENN</i>	0.9560
• SMOTE Tomek	0.9799



Cluster Centroids + SMOTE Tomek are used to:

- undersample *Normal*, *DoS*, *Probe* and *R2L*
- oversample *U2R*

For each 5 preprocessed datasets:

- Decision Tree
- Random Forest
- Extra Tree
- Extra Trees
- k-NN
- SVM
- Logistic regression
- Naive Bayes
- Gradient boosting
- Multilayer perceptron

Parameters
optimization

Tree-structured



Parzen Estimator

Trained classifiers
with optimized
parameters

For each 5 preprocessed datasets:

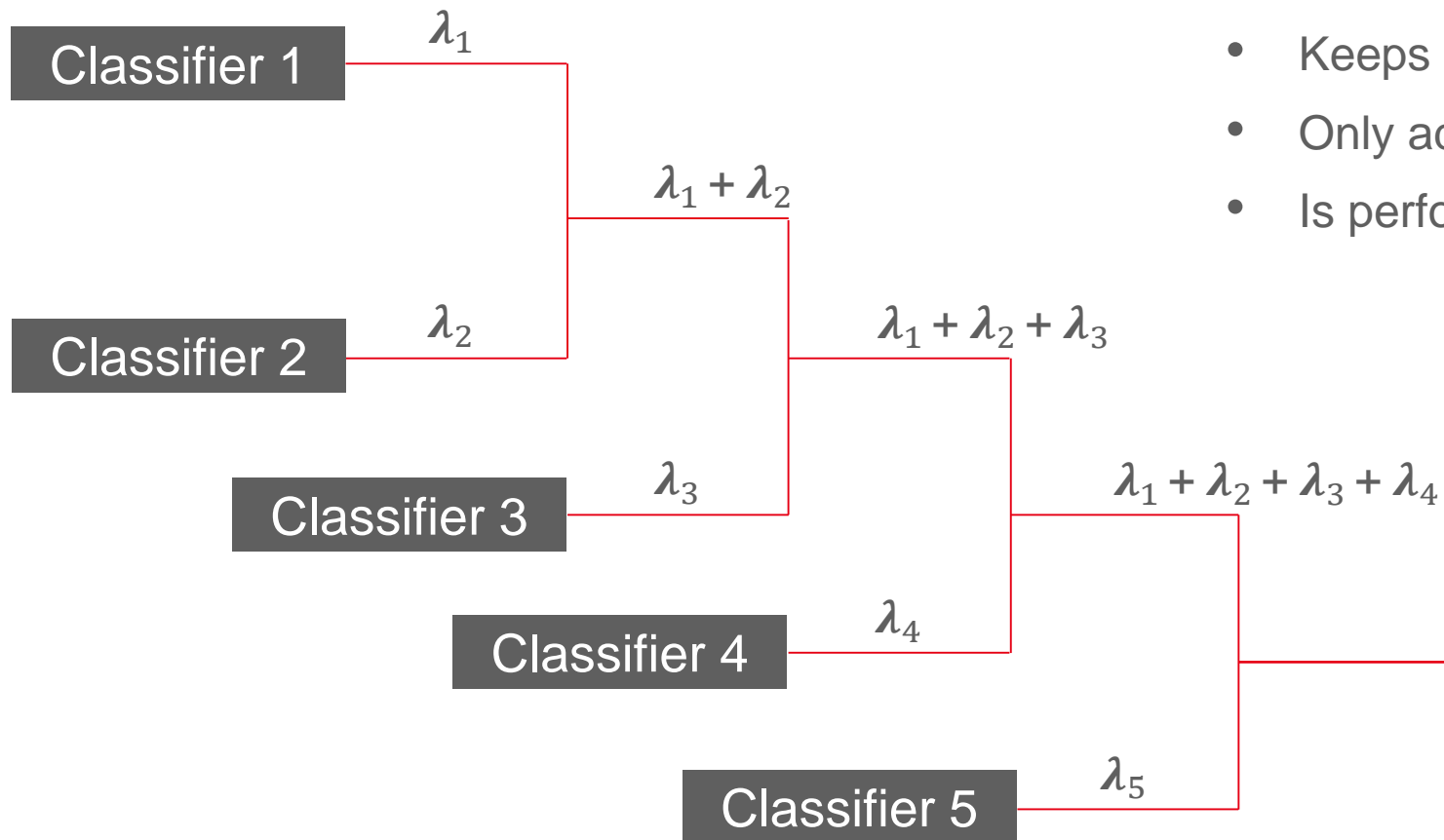
- Decision Tree p_1
- Random Forest p_2
- Extra Tree p_3
- Extra Trees p_4
- k-NN p_5
- SVM p_6
- Logistic regression p_7
- Naive Bayes p_8
- Gradient boosting p_9
- Multilayer perceptron p_{10}

Ensemble

$$p_{class} = \sum_{i=1}^N \lambda_i p_i$$



How to
choose the
best λ_i ?



Each step:

- Tests 2000 values for λ_i (with $\lambda_j = 1 - \lambda_i$)
- Keeps the best λ_i and λ_j in terms of accuracy
- Only adds classifiers that improve accuracy
- Is performed on a validation set

$$p_{class} = \sum_{i=1}^5 \lambda_i p_i$$

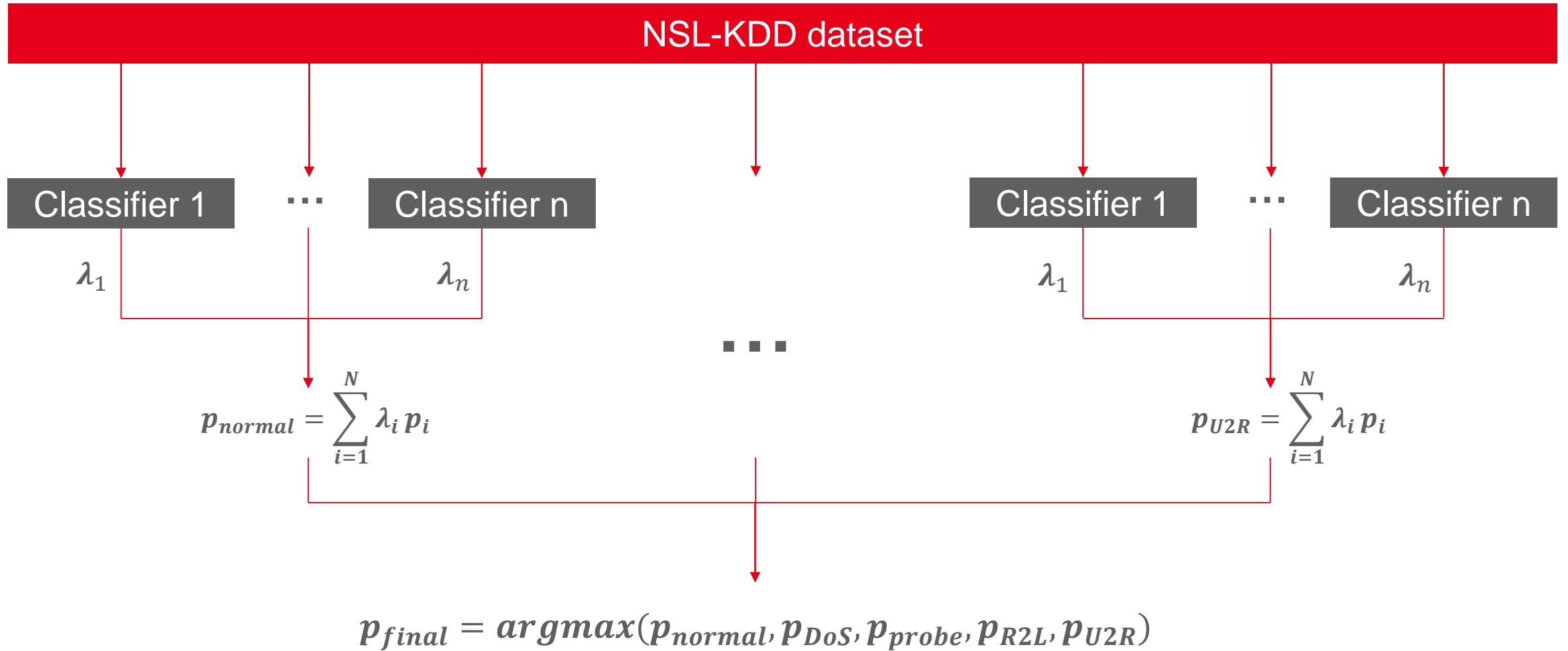
CLASSIFIER	ROC AUC score (test set)	λ_i
MLP3	0.8953	0.29
Decision Tree	0.7446	0.14
Extra Trees	0.8345	0.14
k-NN	0.8178	0.14
MLP1	0.9022	0.14
MLP2	0.8953	0.14

Example:

Probe meta-classifier

$$\begin{aligned}
 p_{probe} = & 0.56 \times p_{MLP1} & + \\
 & 0.05 \times p_{MLP3} & + \\
 & 0.01 \times p_{MLP2} & + \\
 & 0.11 \times p_{k-NN} & + \\
 & 0.14 \times p_{Decision\ Tree} & + \\
 & 0.12 \times p_{Extra\ Trees} &
 \end{aligned}$$

Probe ROC AUC score = 0.9184
(0.9022 at best without ensemble)



NSL-KDD (global accuracy = 86.59%)

	<i>Normal</i>	<i>DoS</i>	<i>Probe</i>	<i>R2L</i>	<i>U2R</i>
Accuracy	95.20%	90.91%	95.87%	91.73%	99.48%
TPR	88.85%	95.99%	86.70%	50.02%	43.28%
FPR	1.66%	12.94%	3.03%	2.15%	0.36%
F1 score	0.9245	0.9009	0.8185	0.6076	0.3295
AUC ROC	0.9359	0.9153	0.9184	0.7394	0.7146

Comparison of results on NSL-KDD

<i>Study</i>	<i>Accuracy</i>	<i>FPR</i>
Our solution	86.59%	1.66%
Two-level classifier ensemble	85.016%	12.6%
Bagging (J48) + feature selection	84.25%	2.79%
GAR-forest + feature selection	85.05%	12.2%
SVM + feature selection	82.37%	15%

- Very high global accuracy and good FPR
- R2L and U2R attacks are still poorly detected
- All optimization steps are fully automated
- Can be applied to:
 - Other datasets (CICIDS 2017)
 - Real-life networks (for classification)

THANK YOU!

Commissariat à l'énergie atomique et aux énergies alternatives
Institut List | CEA SACLAY NANO-INNOV | BAT. 861 – PC142
91191 Gif-sur-Yvette Cedex - FRANCE
www-list.cea.fr

Établissement public à caractère industriel et commercial | RCS Paris B 775 685 019