

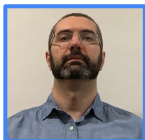


Using Stratified Sampling to Improve LIME Image Explanations

Authors:



UNIVERSITÀ
DI TORINO



Elvio G. Amparore
elviogilberto.amparore@unito.it

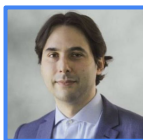


Muhammad Rashid
muhammad.rashid@unito.it

RuleX



Enrico Ferrari
enrico.ferrari@rulex.ai



Damiano Verda
damiano.verda@rulex.ai

Fundings:

EU Horizon-2020 ECSEL-JU project

NEXT
PERCE
PTION



<https://nextperception.eu/>

Topic of the presentation

- Consider linear explanations of **image data**
- Many popular techniques: Grad-CAM, SHAP, LIME, etc...
- We focus on **model-agnostic** systems that generate linear explanations by probing a black-box model using perturbations of the image input.
- Focus on **LIME Image** and its **sampling strategy**.
- LIME Image employs **Monte Carlo sampling** to generate synthetic neighborhoods

How LIME Image works

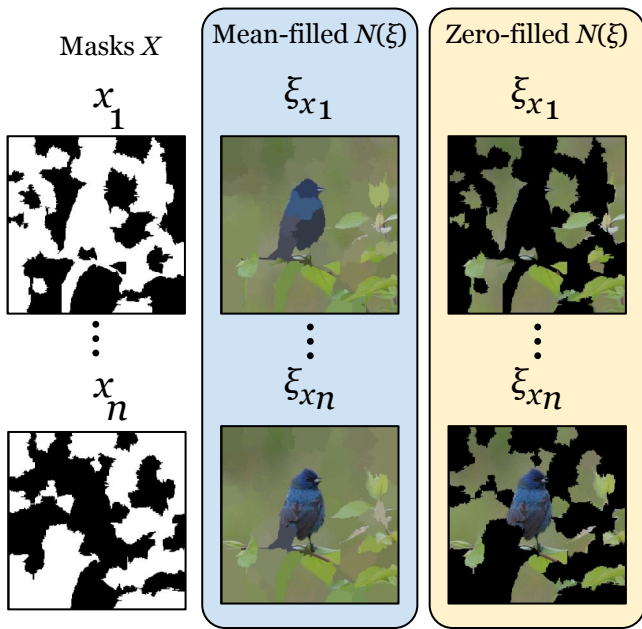
(A) Input image ξ



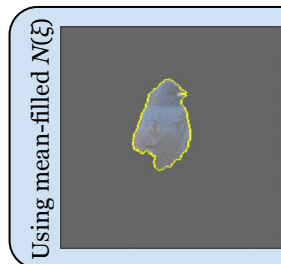
(B) Superpixels
(interpretable
representation)
 $k=84$



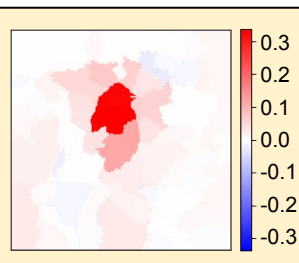
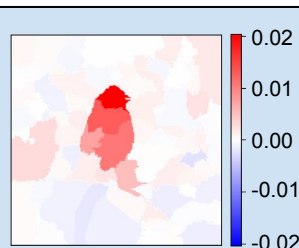
(C) Neighborhood generation
from the set of n masks X



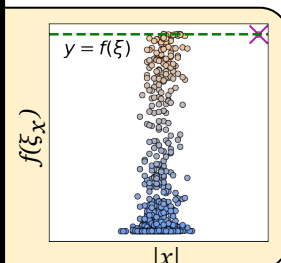
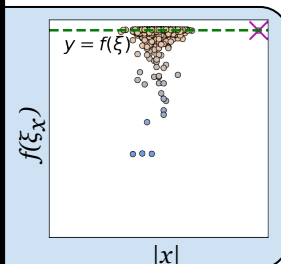
(D) LIME Image
explanations



(E) Superpixel
importances β



(F) Dependent
variable
distribution



Synthetic neighborhood generation

Initial image $\xi \rightarrow$ divided into k superpixels

Superpixel masking: $x \in \{0, 1\}^k$ generates a perturbed image ξ_x

In LIME Image, masks are sampled using an **unbiased Monte Carlo** strategy:

$$x[i] \sim B(0.5), \quad 1 \leq i \leq k$$

where $B(p)$ is a Bernoulli-distributed random variable having probability $p = 0.5$

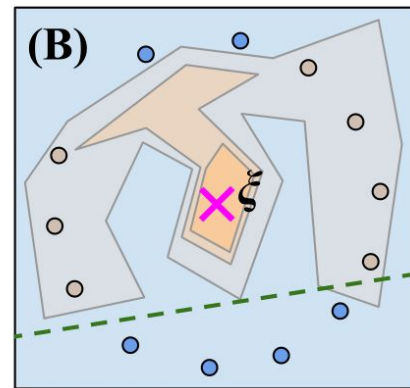
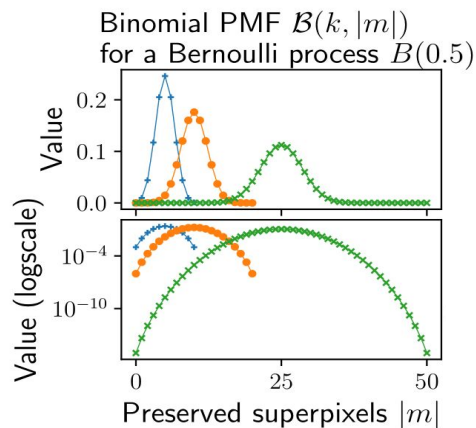
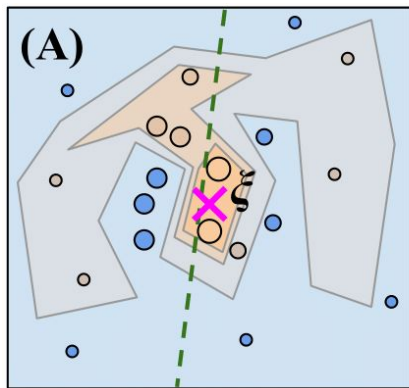
$N(\xi)$ = synthetic neighborhood of n perturbed images

Dependent variables: $Y = \{f(\xi_x) \mid \xi_x \in N(\xi)\}$

black box model being explained

The problem with the unbiased Bernoulli distribution

- Ideally, the synthetic neighborhood $N(\xi)$ should provide a “good enough” coverage of the variations around x .



- May result in **under-representation** of the neighborhood

* Image freely inspired by: Ribeiro, M. T.; Singh, S.; and Guestrin, C. 2016. Why should I trust you? Explaining the predictions of any classifier. In Proceedings of the 22nd ACM SIGKDD int. Conf. 1135–1144.



(A) Input image ξ

Class: *hyena*

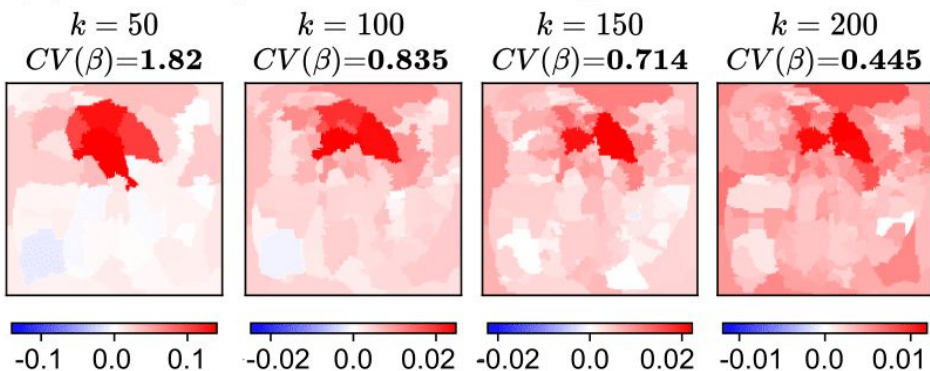
Probability: 99.46%

num_samples n = 1000

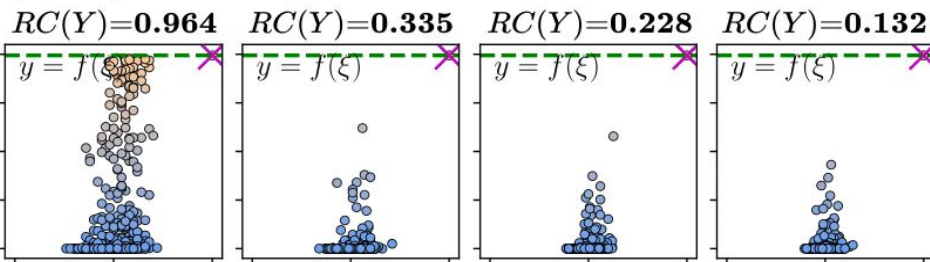
Using mean-filled $N(\xi)$

| k | max_dist |
|-----|-------------|
| 50 | 8.691 |
| 100 | 4.956 |
| 150 | 4.092 |
| 200 | 3.632 |

(B) Feature importances β for four segmentations.



(C) Dependent variable distributions.



Dependent variable undersampling results in confused explanations

Coefficient of Variation of the explanation β

$$CV(\beta) = \frac{\sigma_{\beta}}{\mu_{\beta}}$$

Range Coverage of the values of the dependent variable Y in the synthetic neighborhood.

$$RC(Y) = \frac{IQR_{1-99}(Y)}{f(\xi)}$$

- Almost no sample is close to ξ .
- Samples drawn from $B(p)$ have all about 50% of the superpixels masked.
- Under-representation of the local behaviour of the black-box model f .

Sampling relevance

Bernoulli distribution is not the only option.

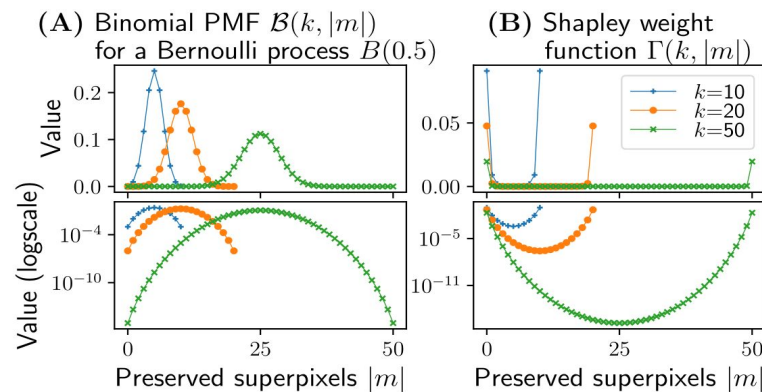
Shapley theory uses a different distribution:

Shapley Importance
function

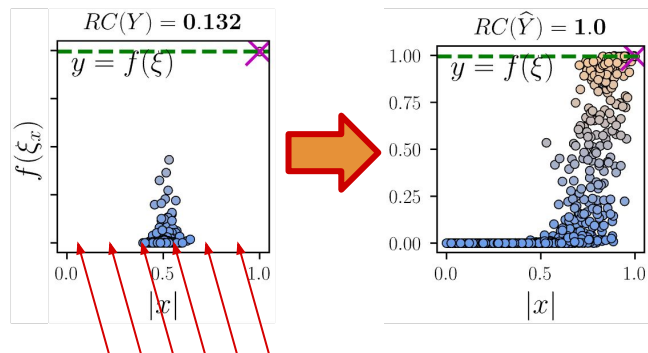
$$\Gamma(k, |x|) = \frac{1}{(k+1) \binom{k}{|x|}}$$

Shapley Importance is
the reciprocal of the
Bernoulli distribution

$$\mathcal{B}(k, |x|) \cdot \Gamma(k, |x|) = \frac{\binom{k}{|x|} p^{|x|} (1-p)^{k-|x|}}{(k+1) \binom{k}{|x|}} = \frac{0.5^k}{k+1}$$



Proposed methodology: stratified sampling



The goal is to draw samples to cover the entire $|x|$ sampling space

- Consider a stratified partitioning
- $\mathcal{X}^{(i)}$ = set of possible masks having $|x| = i$
- Stratum i size is known a priori:

$$|\mathcal{X}^{(i)}| = \binom{k}{i}, \quad 0 \leq i \leq k$$

- Oversampled probability

$$Prob\{x \in \mathcal{X}^{(i)} \mid x \in \hat{X}\} = \frac{1}{k+1}$$

- Adjustment factors

$$adj(i) = \frac{Prob\{x \in \mathcal{X}^{(i)} \mid x \in X\}}{Prob\{x \in \mathcal{X}^{(i)} \mid x \in \hat{X}\}} = \frac{(k+1)\binom{k}{i}}{2^k}$$

Monte Carlo vs. Stratified Sampling

LIME Image uses a **simple linear homoscedastic model**

$$Y = X \cdot \beta + \epsilon$$

The explanation coefficients β results from

$$\beta = (X^T W X)^{-1} X^T W Y$$

β coefficients may vary by stratum.

We adopt instead a **mixture model**

$$\hat{Y}^{(i)} = \hat{X}^{(i)} \cdot \hat{\beta}^{(i)} + \hat{\epsilon}^{(i)}$$

The explanation coefficients β results from

$$\hat{\beta} = (\hat{X}^T \hat{W} \hat{X})^{-1} \hat{X}^T \hat{W} \hat{Y}$$

where \hat{W} accounts for the adjustment factors that correct the bias introduced by oversampling the distribution tails.

Impact of stratified sampling in LIME Image

Case (A): The mean and variance of $\hat{\beta}^{(i)}$ are independent from the strata.

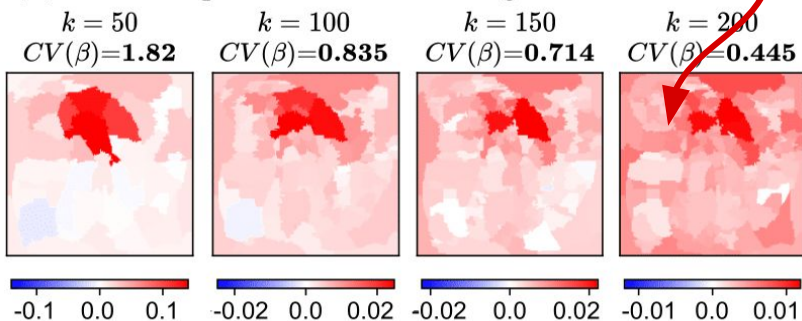
- ⇒ weighted regression model is not needed.
Monte Carlo and stratified sampling should behave similarly.
Unlikely to happen using complex black-box machine learning models.

Case (B): The mean and variance of $\hat{\beta}^{(i)}$ varies by stratum.

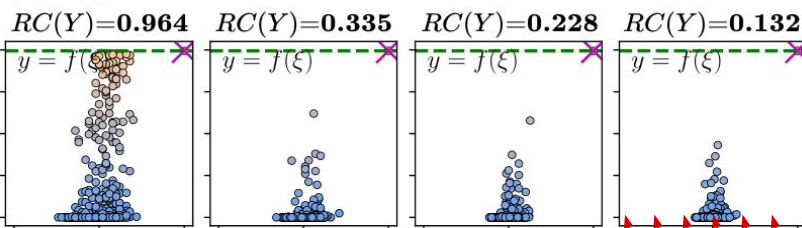
- ⇒ weighted regression model is highly advisable (DuMouchel and Duncan 1983).
Monte Carlo will perform badly, stratified sampling is relevant.
Common scenario for complex models and/or large number of superpixels.

Monte Carlo sampling

(B) Feature importances β for four segmentations.



(C) Dependent variable distributions.

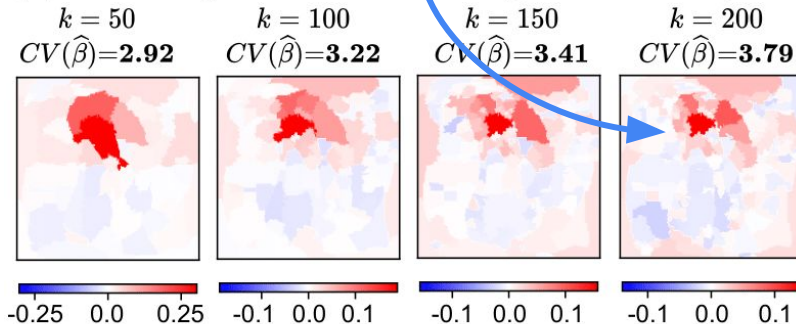


undersampling

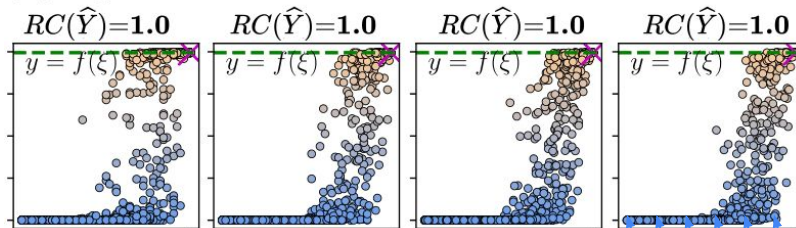
Stratified sampling

feature attribution
is meaningful

(B) Feature importances $\hat{\beta}$ for four segmentations of *hyena*.



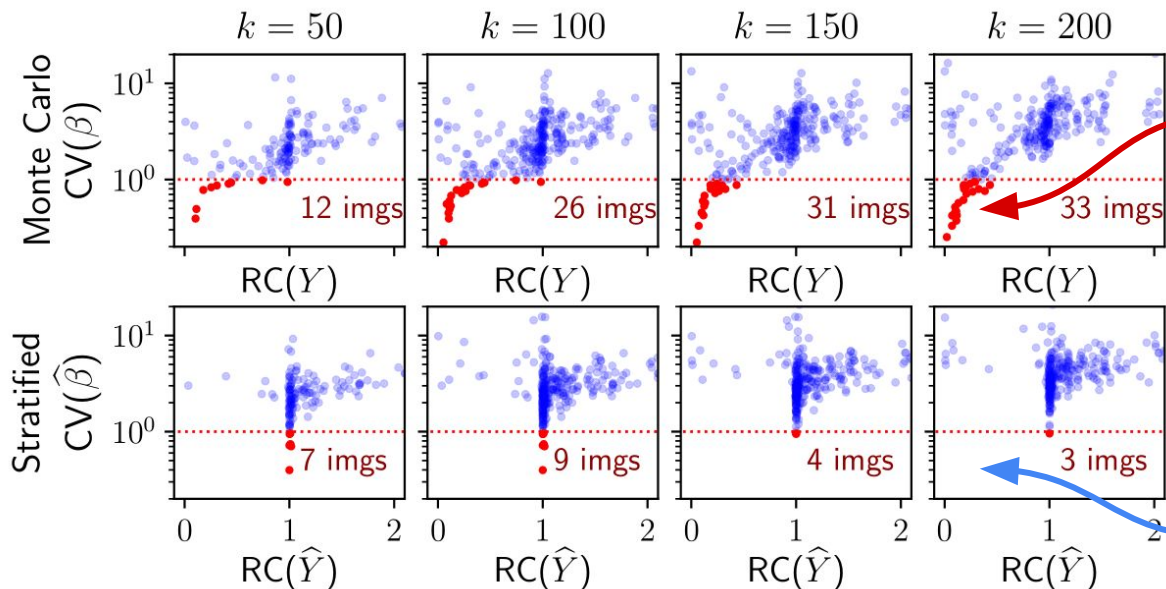
(C) Dependent variable distributions.



samples from
all strata








Evaluation on the *ImageNet Object Localization* Dataset

150 images; four different choices of superpixels ($k = 50, 100, 150, 200$)
n=1000 samples (average of 10 runs). Using ResNet-50 model.



About 1 image out of 5 suffers from severe undersampling using the default Monte Carlo sampling of LIME-Image

Misbehaviours are corrected using stratified sampling

| | | Monte Carlo sampling | | | | Stratified sampling | | | |
|----------------|--|----------------------|---------|-------------|---------|---------------------|---------------|-------------------|---------------|
| | | (A) | | (B) | | (C) | | (D) | |
| Image | | $k = 50$ | | $k = 200$ | | $k = 50$ | | $k = 200$ | |
| | | $CV(\beta)$ | $RC(Y)$ | $CV(\beta)$ | $RC(Y)$ | $CV(\hat{\beta})$ | $RC(\hat{Y})$ | $CV(\hat{\beta})$ | $RC(\hat{Y})$ |
| Low CV Values | orange  71.3% | 0.387 | 0.0833 | 0.178 | 0.0217 | 2.12 | 1.0 | 3.54 | 1.01 |
| | wardrobe  31.4% | 0.135 | 0.0182 | 0.164 | 0.00908 | 1.72 | 1.0 | 3.04 | 0.993 |
| | milk can  76.2% | 0.967 | 0.357 | 0.571 | 0.145 | 3.04 | 1.14 | 4.71 | 1.14 |
| | lynx  50.4% | 1.57 | 1.19 | 0.464 | 0.137 | 3.05 | 1.71 | 4.79 | 1.58 |
| | ringneck snake  36.2% | 1.10 | 0.365 | 0.221 | 0.0238 | 3.45 | 1.3 | 5.11 | 1.19 |
| High CV Values | chickadee  99.8% | 4.05 | 0.999 | 5.73 | 0.996 | 3.74 | 1.0 | 4.05 | 1.0 |
| | polecat  46.9% | 7.07 | 1.86 | 5.76 | 1.51 | 6.54 | 1.79 | 5.76 | 1.46 |

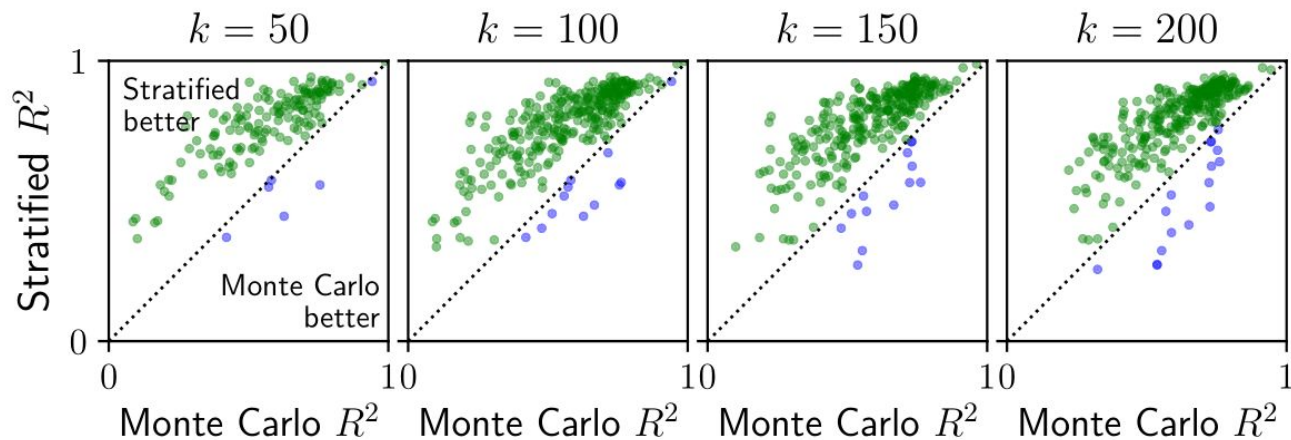
Evaluation on the ImageNet Object Localization Dataset

Inspection shows that poor sampling of synthetic neighborhoods using unbiased Monte Carlo **are not rare**.

Unsurprisingly, explanations built using these (poor) neighborhoods do not identify the relevant portions of the image

Evaluation on the *ImageNet Object Localization Dataset*

In general the dependent variable sampled using Stratified Sampling is a better explainer than the one sampled using Monte Carlo (tested using average R^2 coefficients of the linear regressors built by LIME).



Conclusions

- Reformulation of LIME Image sampling strategy (not restricted to image data) for stratified sampling.
- Drawing lessons from the Shapley theory.
- Empirical evaluation shows that Monte Carlo undersampling is not rare, and stratified sampling provides practical improvements, at no additional cost.

Possible improvements

- Consider regularization factors for ridge regression.
- Mixed model could be improved using uniform weights for the strata.

Code Availability

- https://github.com/rashidrao-pk/lime_stratified
- <https://github.com/rashidrao-pk/lime-stratified-examples>





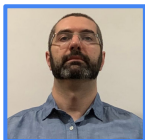
Using Stratified Sampling to Improve LIME Image Explanations

Thank you

Authors:



UNIVERSITÀ
DI TORINO



Elvio G. Amparore
elviogilberto.amparore@unito.it

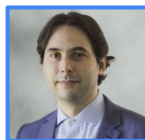


Muhammad Rashid
muhammad.rashid@unito.it

RuleX



Enrico Ferrari
enrico.ferrari@rulex.ai



Damiano Verda
damiano.verda@rulex.ai

Fundings:

EU Horizon-2020 ECSEL-JU project

NEXT
PERCE
PTION



<https://nextperception.eu/>