# OCTET: Object-aware counterfactual explanations

Change the model's decision

Decision model $M$ → *Left* ✅

Query $x^q$

Small and meaningful changes

Decision model $M$ → *Left* ❌
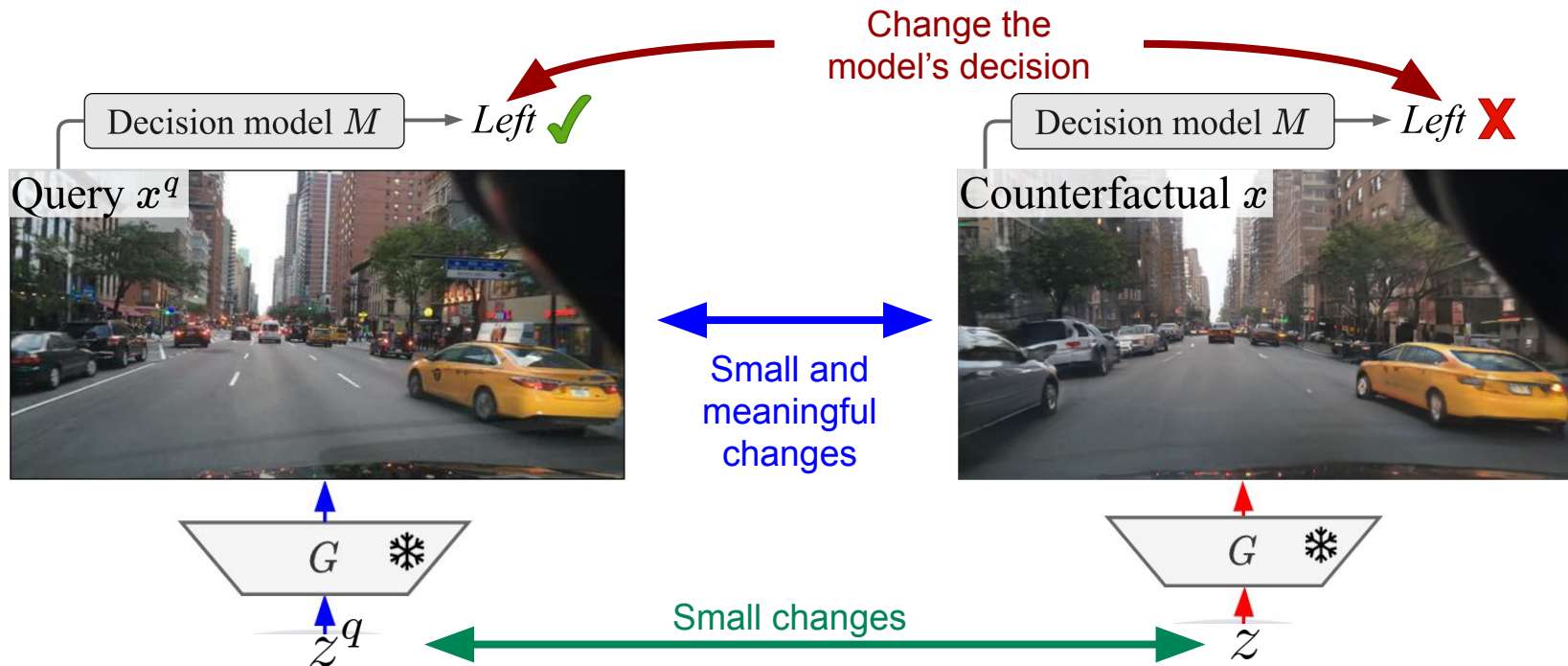
Counterfactual $x$

**Contributions:**

New object-aware framework that seeks counterfactual explanations for **compositional images and complex scenes**

Using an object-centric latent space empowers users to **explore the role of specific objects** in decision models

**User-study** showing usefulness of counterfactual explanations to better anticipate failures and find biases.

Explanations of **semantic segmentation models**

# Counterfactual explanations for <u>image</u> classification models?



$$\mathrm{argmin}_z \mathscr{L}_{\mathrm{distance}}(z^q, z) + \lambda \mathscr{L}_{\mathrm{decision}}(M(G(z)), y)$$

# Dealing with complex compositional scenes

**Style**

| | |
|---|---|
| Color of traffic light | Visibility conditions |
| Road markings | Car rear light |

**Structure**

| | |
|---|---|
| Traffic signs | Car obstacles |
| Pedestrians | … |

**Driving action**

**BlobGan** [Epstein *et al.* ECCV'22]



Blob appearance features

$$\psi$$

$$\phi = c_x, c_y, s, a, \theta$$

Blob spatial parameters

Layout network

$$z \sim \mathcal{N}(0, I)$$

$$\frac{s}{\sqrt{a}} \quad s\sqrt{a}$$
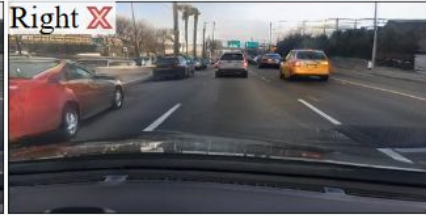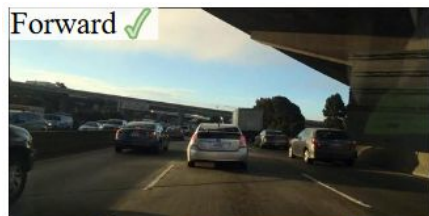
$$\theta$$

$$x$$

$$G$$

# OCTET: overview



① Optimizing blob parameters to get a good image inversion

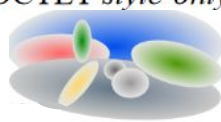② Optimize the blob parameters to change the model decision

# Qualitative results



Forward ✓ | Forward ✗ | Forward ✗ | Forward ✗
Stop ✓ | Stop ✗ | Stop ✗ | Stop ✗
Left ✗ | Left ✓ | Left ✓ | Left ✓
Right ✓ | Right ✗ | Right ✗ | Right ✗
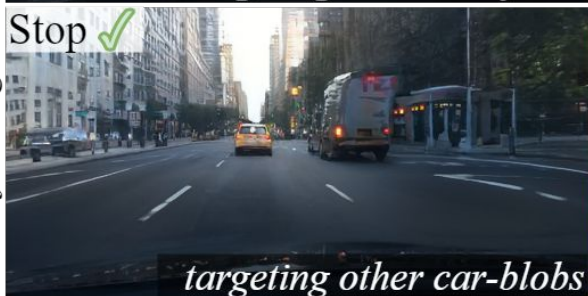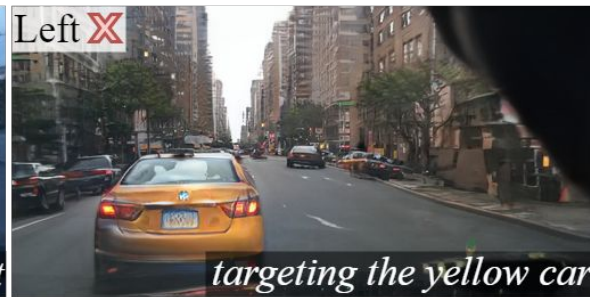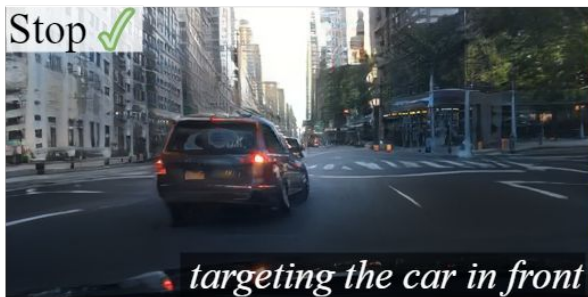
Query | OCTET | OCTET *style-only* | STEEX

ECCV'22

# Targeted counterfactual explanations

Empowering users to explore the role of specific objects in decision models
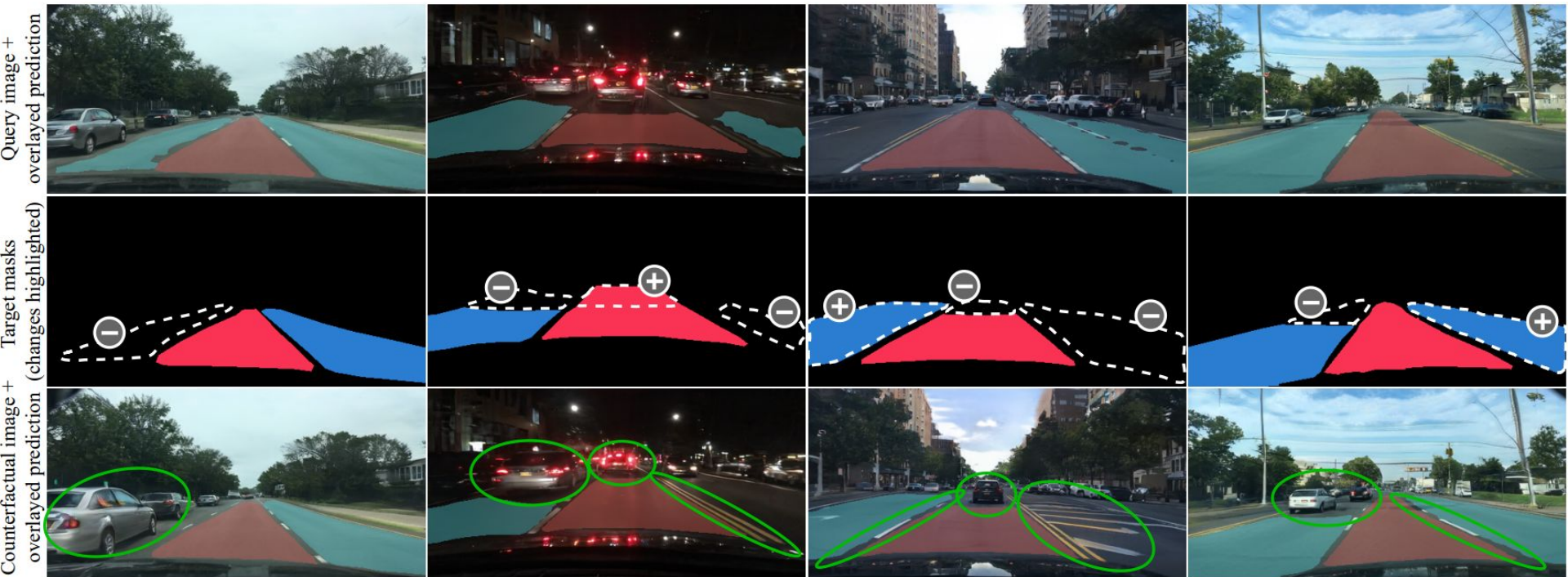
# Explaining a semantic segmentation model

Explaining a DeepLabv3 trained to segment primary and secondary free-spaces
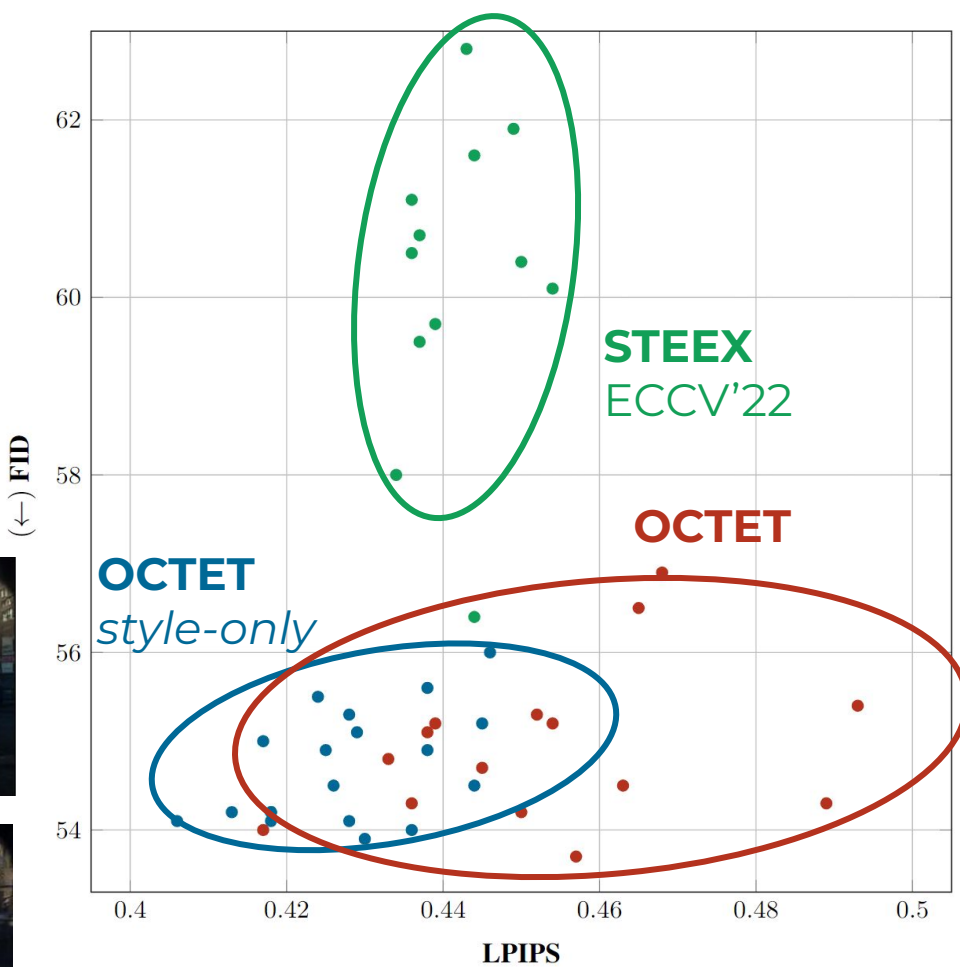
# Quantitative metrics

Are counterfactuals realistic?

→ Fréchet Inception Distance (**FID** (↓))

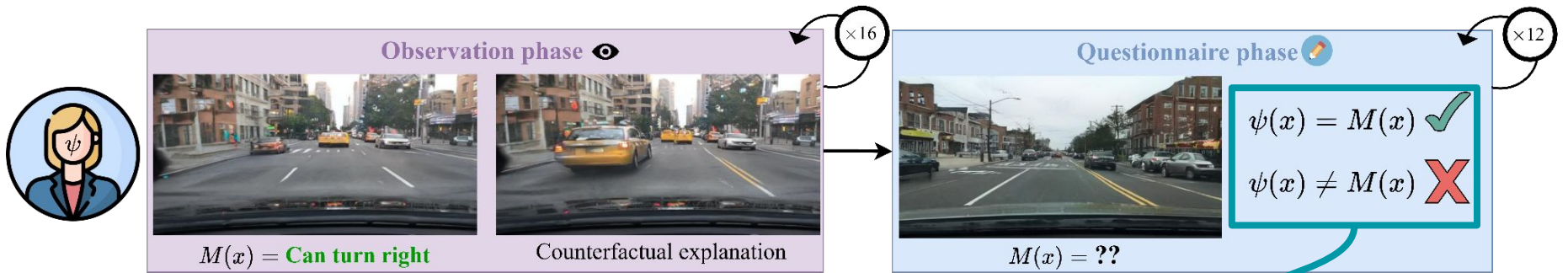Are counterfactuals close to the query image?

→ **LPIPS** distance



Query

OCTET

OCTET *style-only*

STEEX ECCV'22

$$\mathrm{argmin}_x \mathscr{L}_{\text{distance}}(x^q, x) + \lambda \mathscr{L}_{\text{decision}}(M(x), y)$$

# Can counterfactuals help to better "understand" a model?

"Understand" := Ability to predict model's decision on new instances (simulatability)



**Observation phase** 👁

$M(x) =$ **Can turn right**

Counterfactual explanation

**Questionnaire phase** ✏

$M(x) = ??$

$\psi(x) = M(x)$ ✔

$\psi(x) \neq M(x)$ ✘

|  | Cohort size | Replication | Bias Detection |
|---|---|---|---|
| Group with **OCTET** explanations | 20 | **70%** | **65%** |
| Control group **without explanations** during training | 20 | 52.5% | 0% |

p-value = 0.0028

Unknown to the participants, the classifier is flawed: obstacles on both sides of the road influence the "Can turn right" prediction. **Did users find out?**

What I cannot predict, I do not understand: A human centered evaluation framework for explainability method. Fel *et al.*, NeurIPS 2022

# Conclusion

**OCTET is an object-centric framework able to explain classifiers operating on complex and compositional scenes**

github.com/valeoai/octet

WED-PM-258