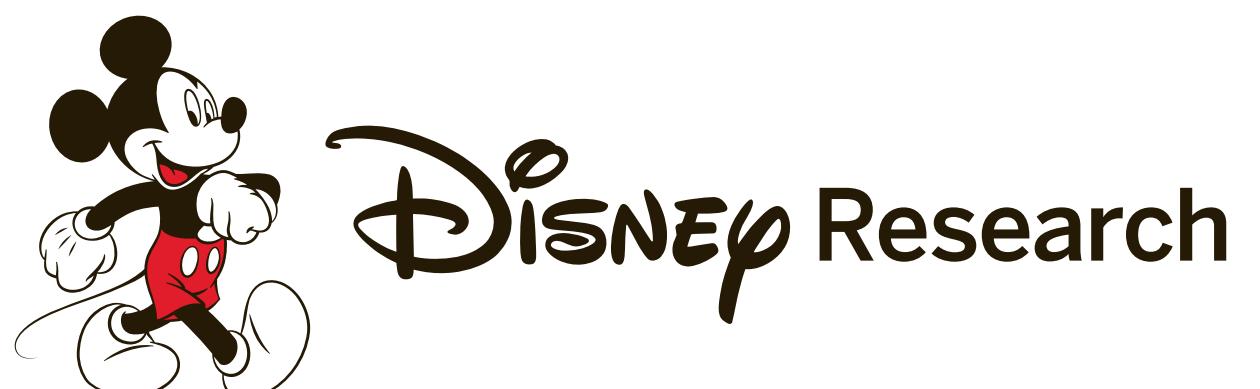




# Semi-supervised Vocabulary- informed Learning

**Yanwei Fu**

**Leonid Sigal**



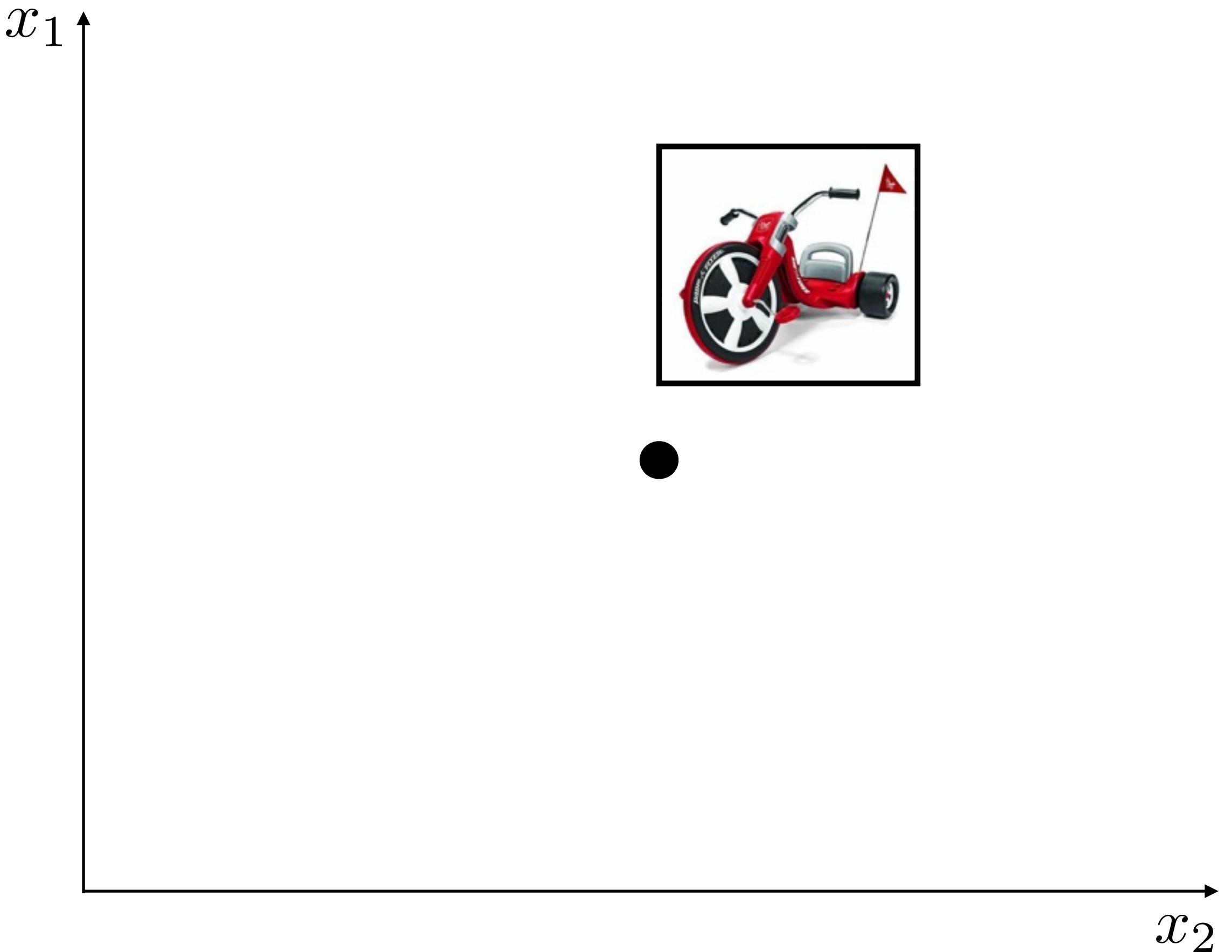
# Supervised Learning

Problem Definition

# Supervised Learning

Problem Definition

Visual feature space



# Supervised Learning

Problem Definition

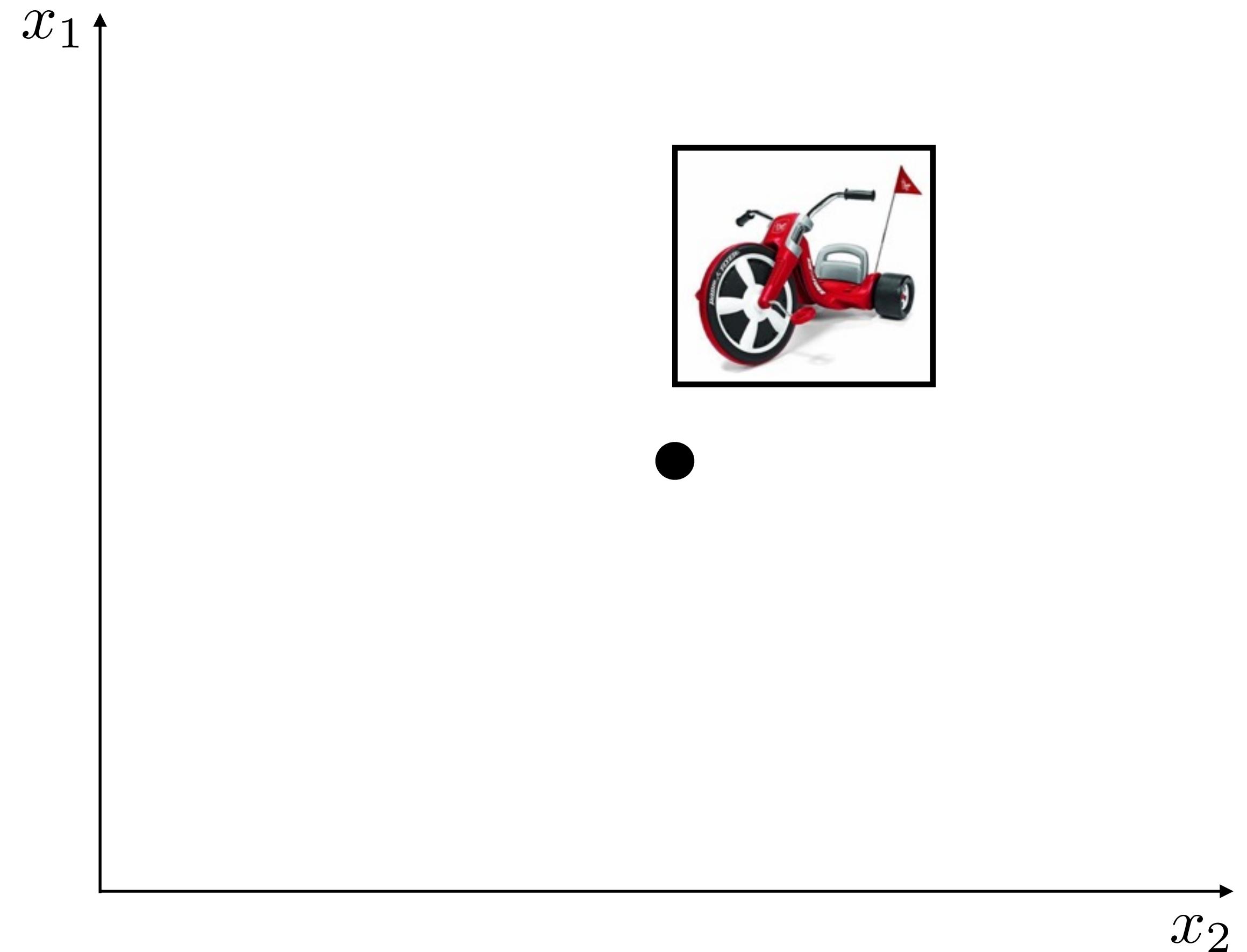


**airplane**

**car**

**unicycle**

**tricycle**



# Supervised Learning

Learning

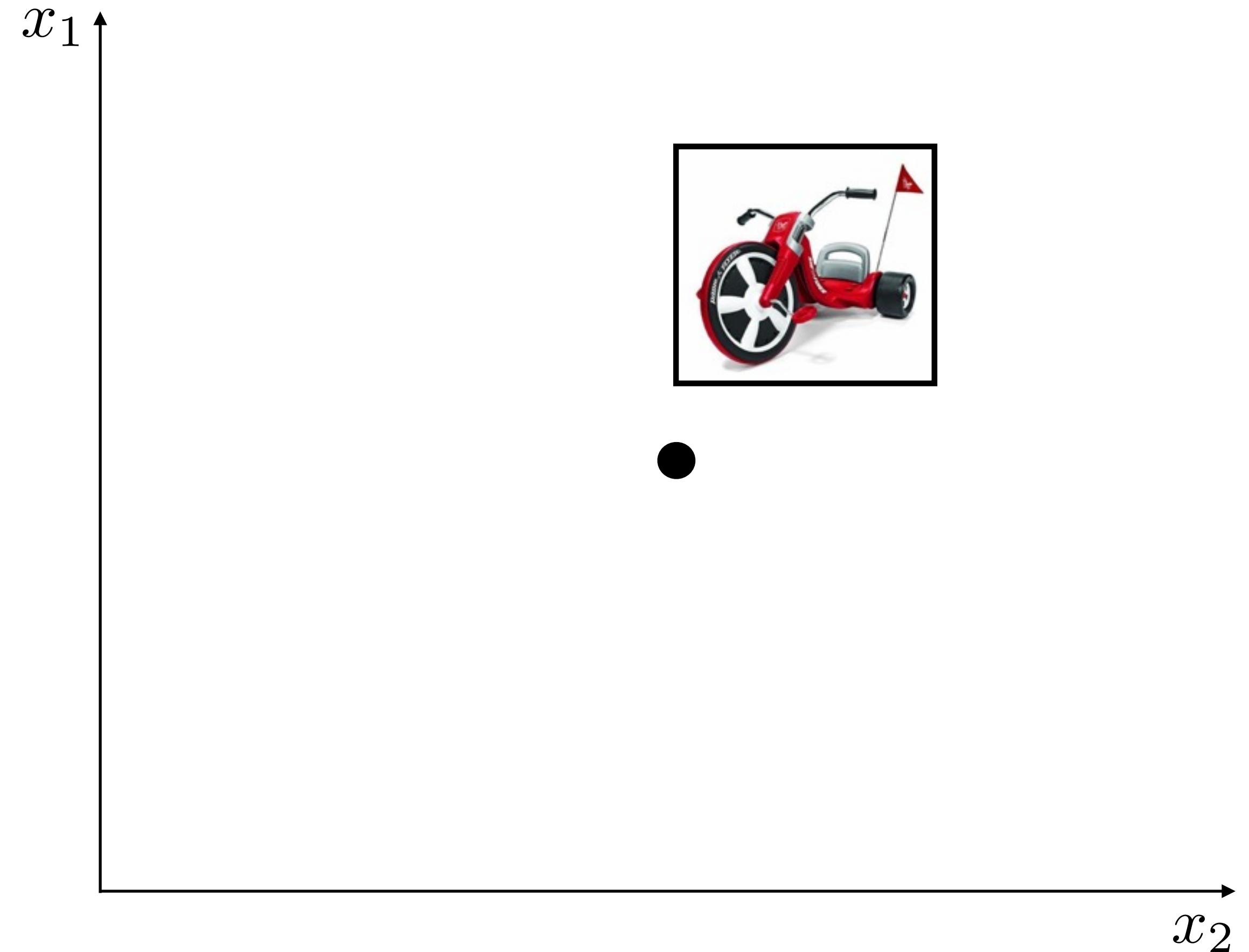


**airplane**

**car**

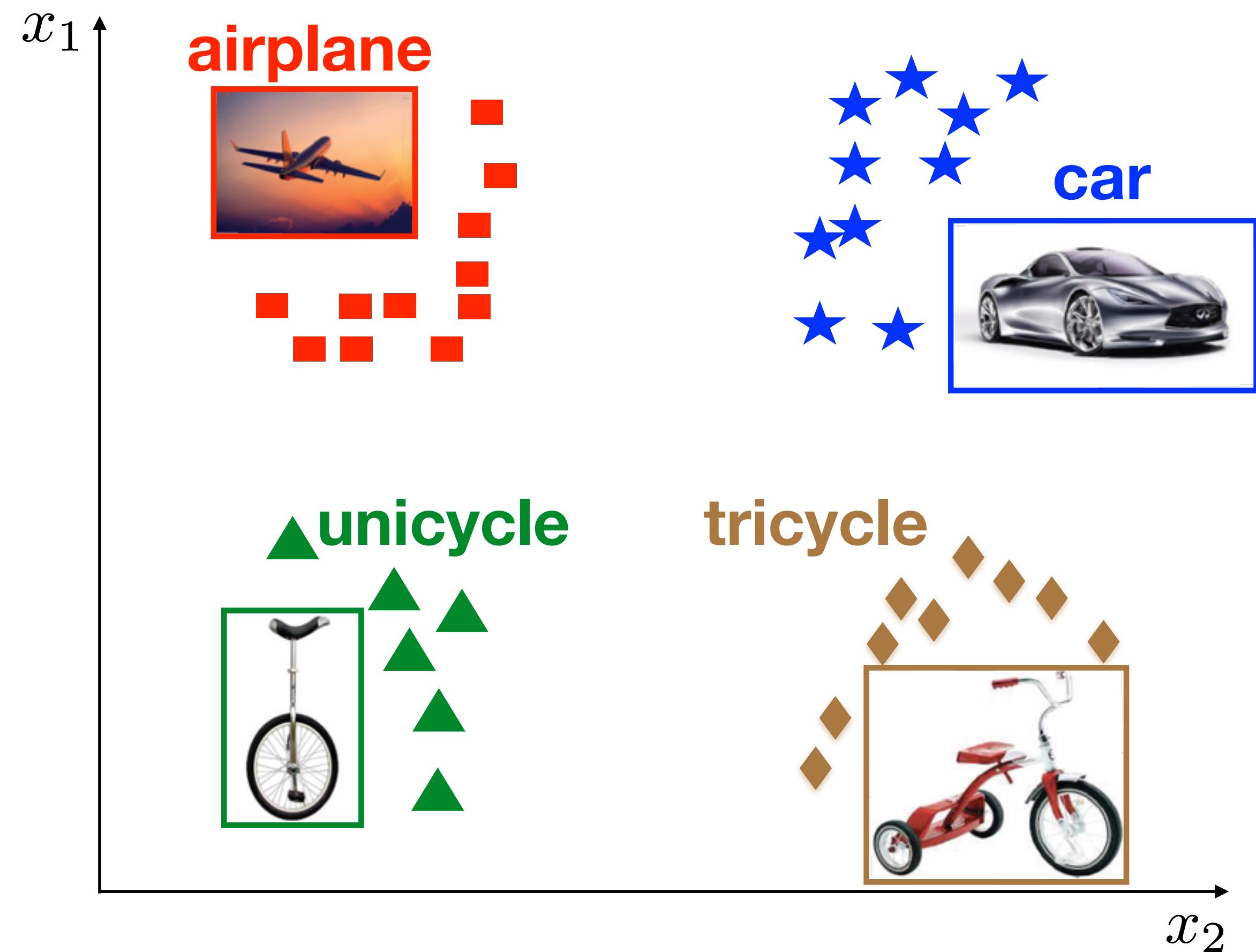
**unicycle**

**tricycle**



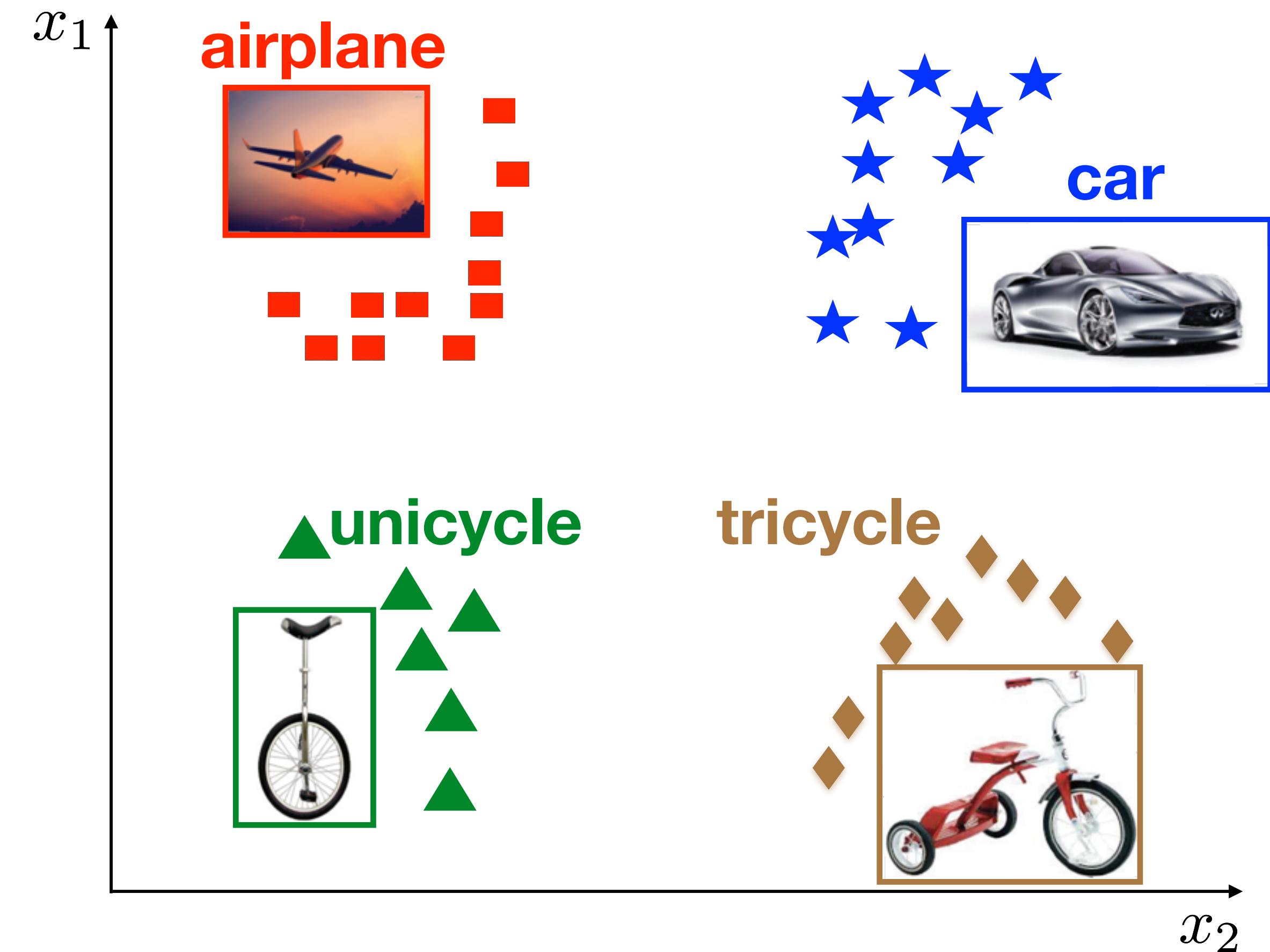
# Supervised Learning

Learning



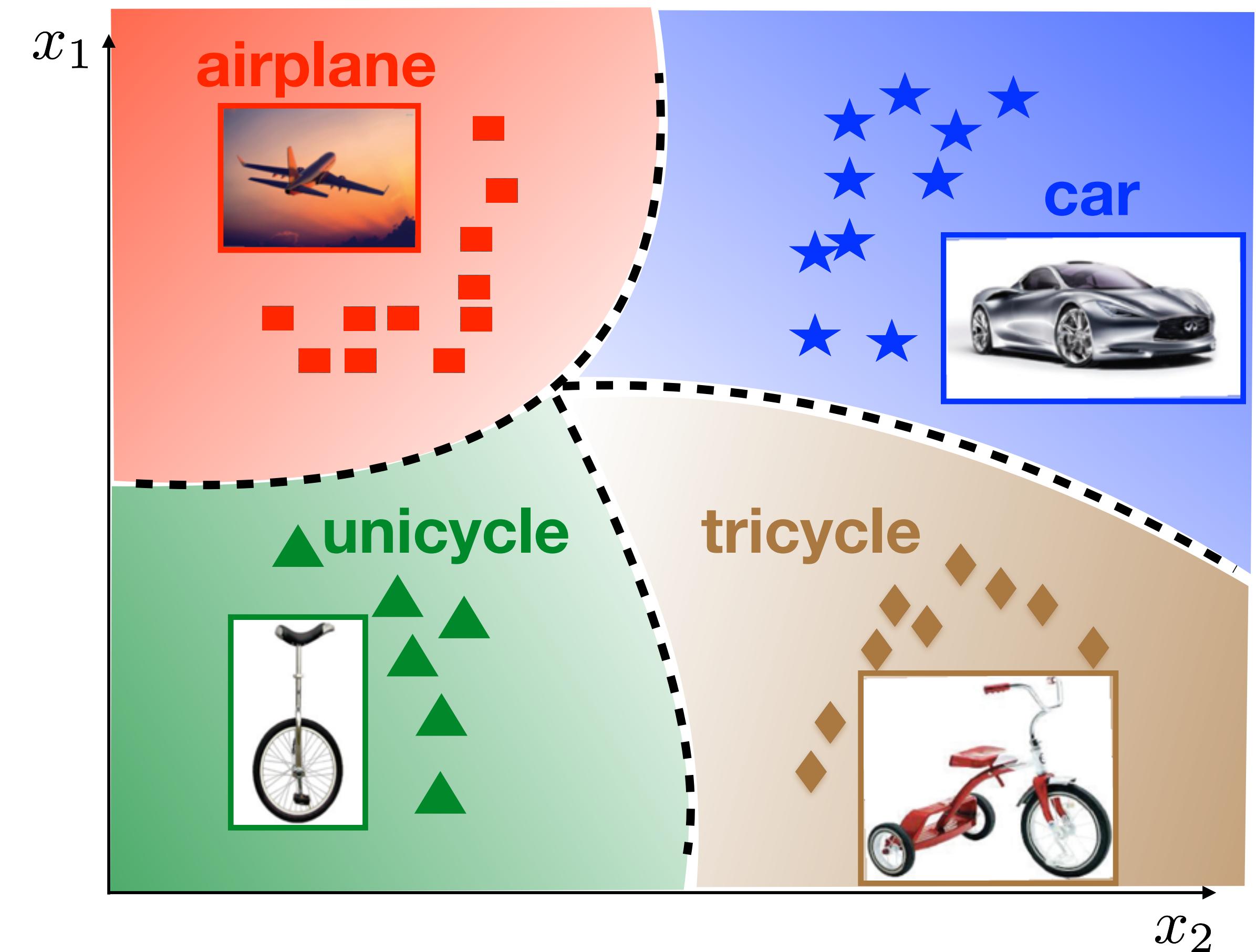
# Supervised Learning

Learning



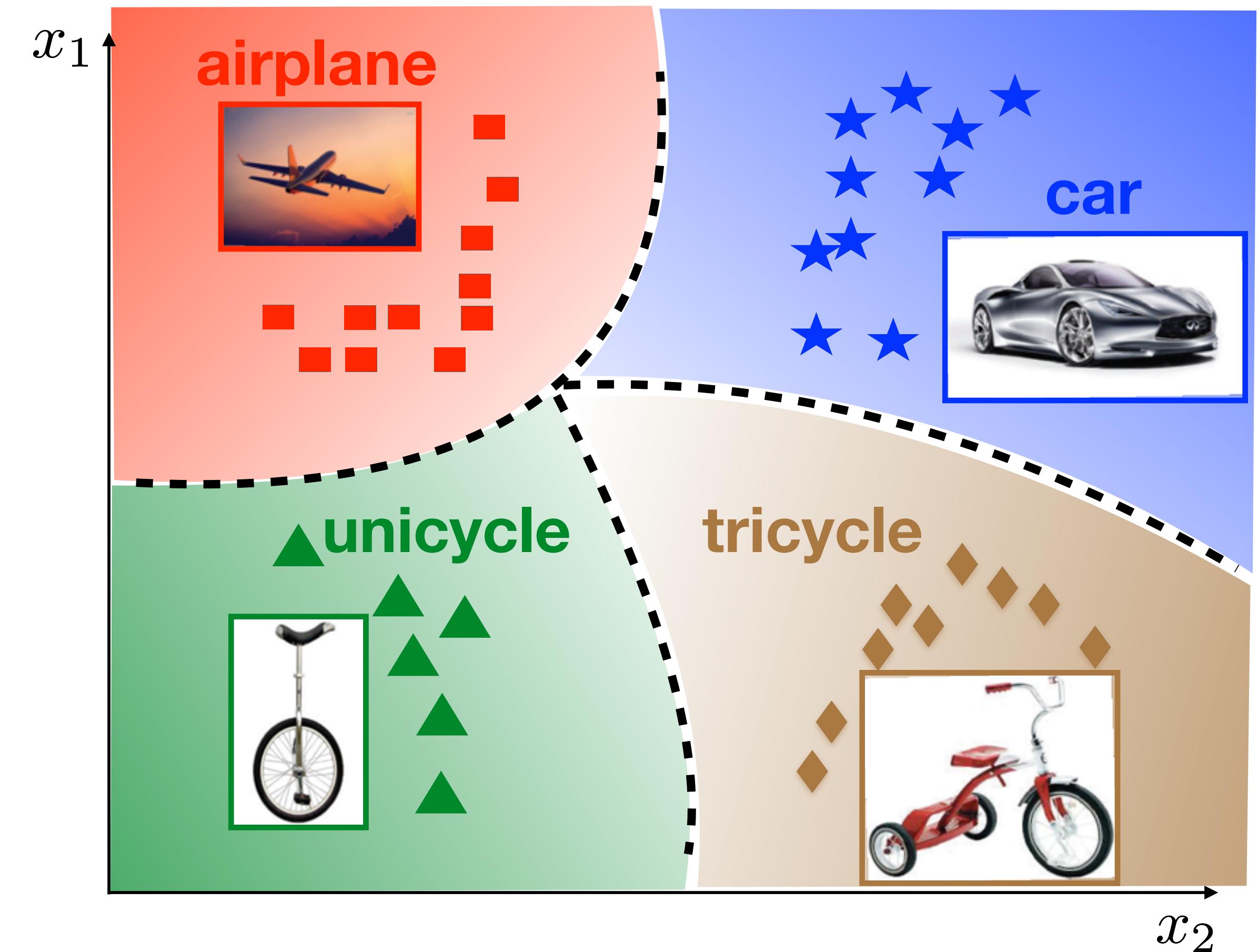
# Supervised Learning

Learning



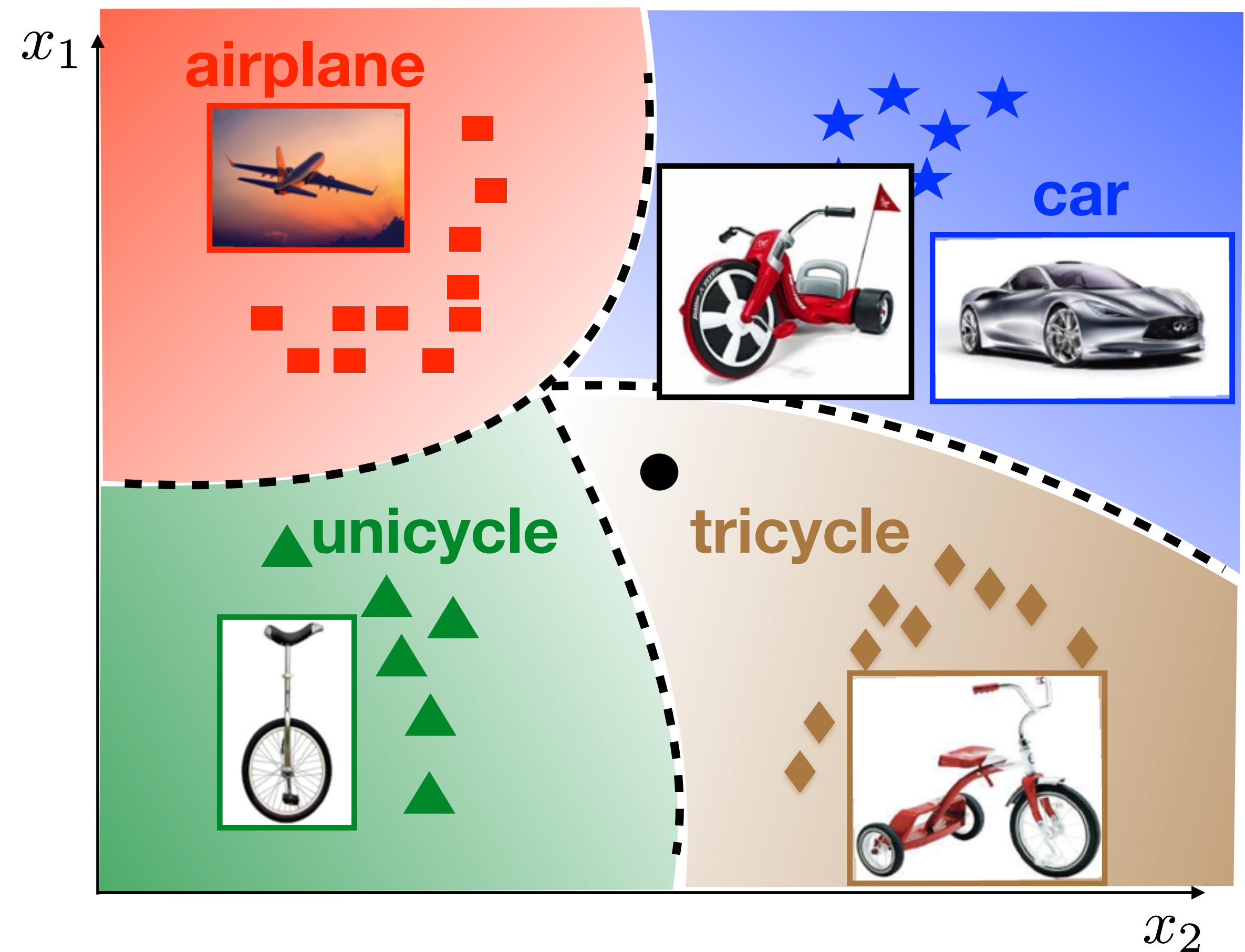
# Supervised Learning

Inference



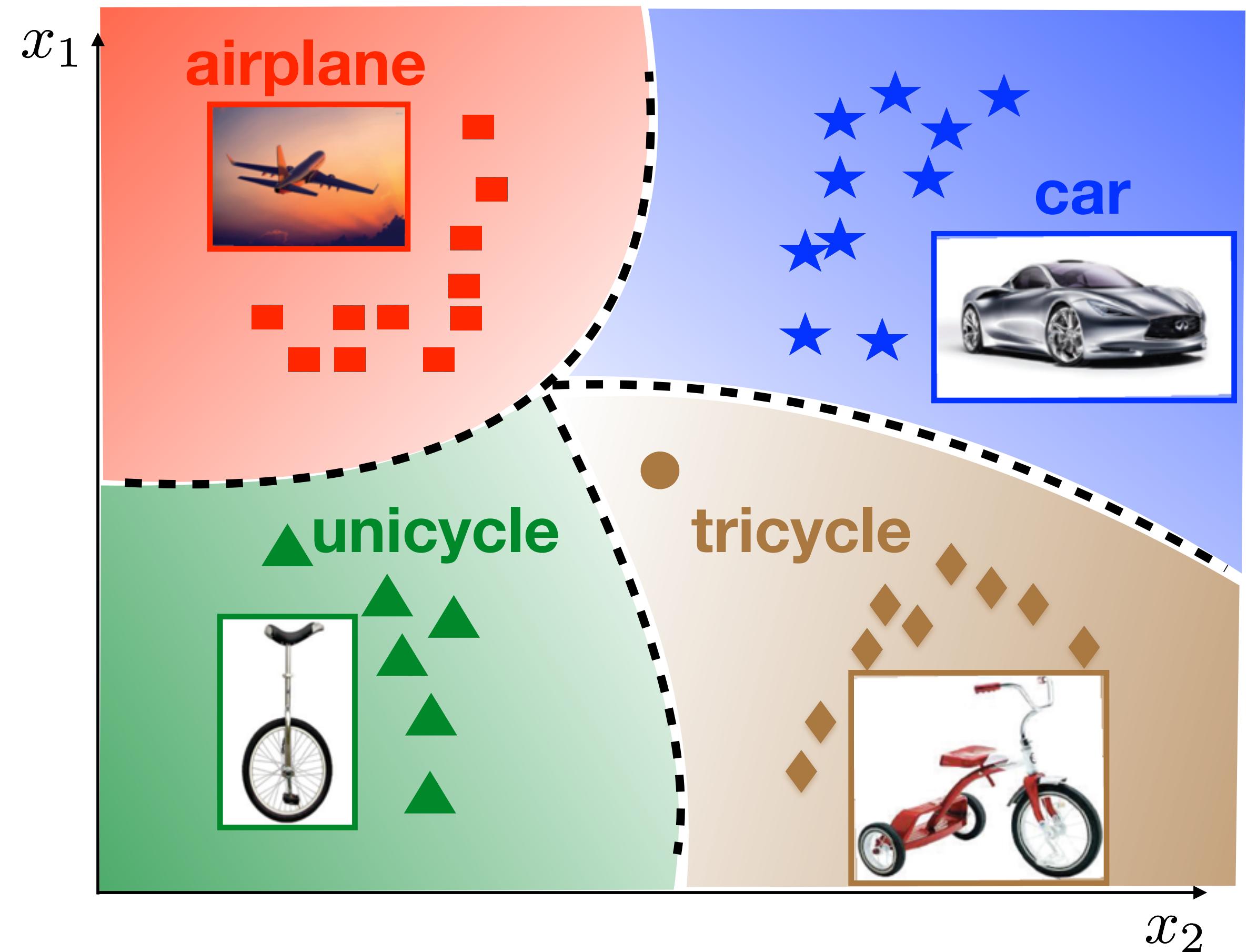
# Supervised Learning

Inference



# Supervised Learning

Inference



# Zero-shot Learning

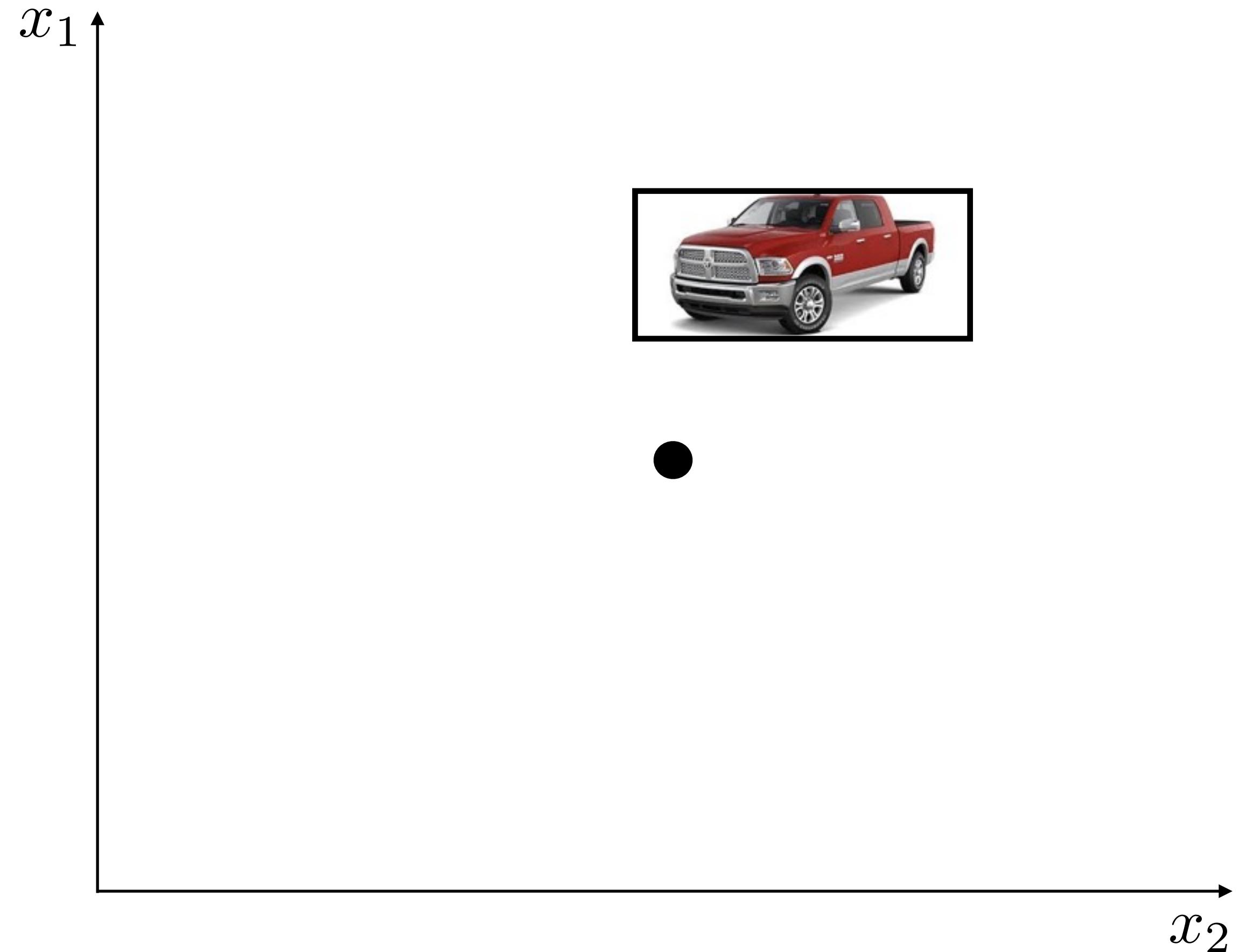
Problem Definition



We do not have any visually labeled instances of what these look like

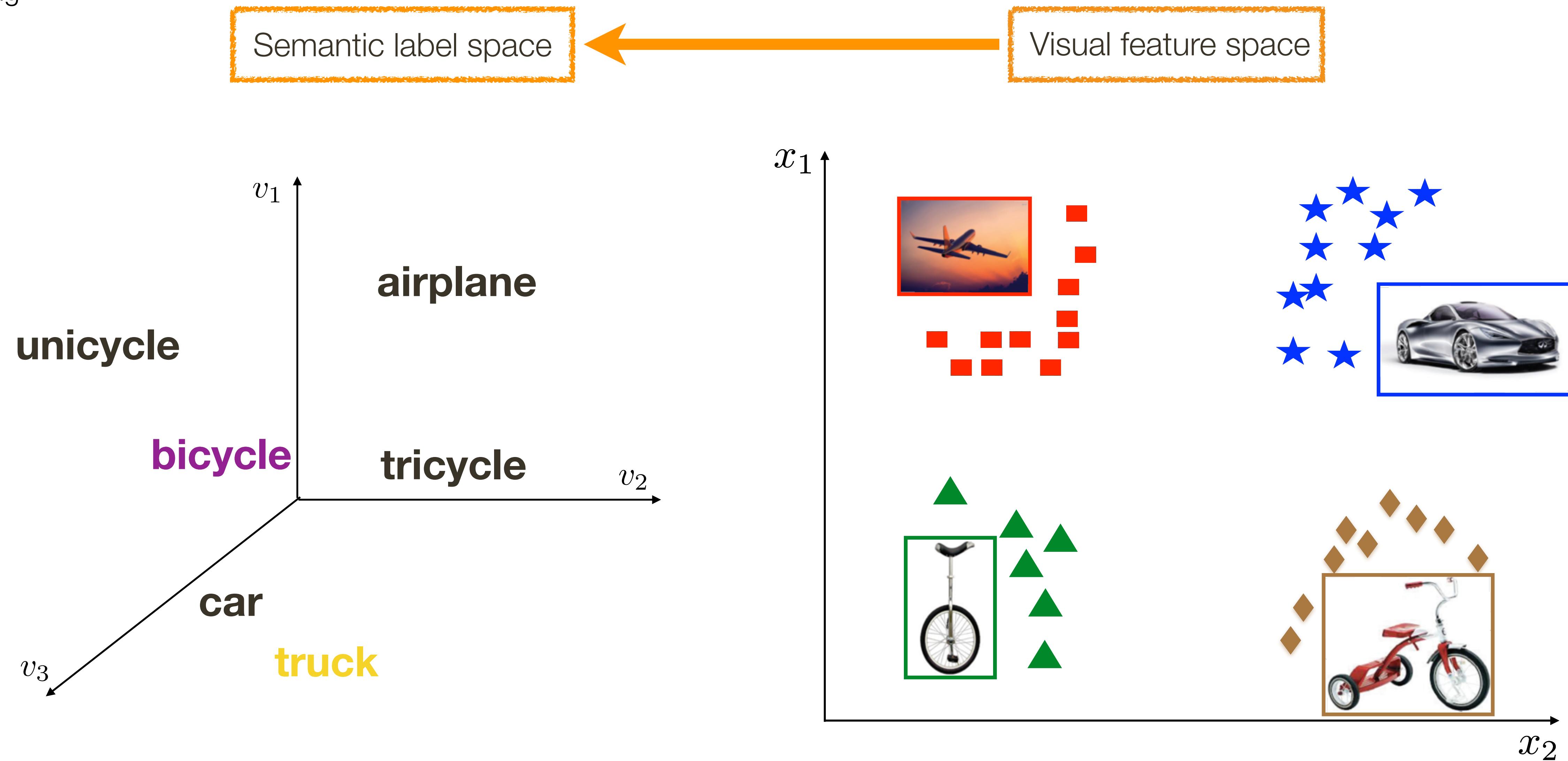
bicycle

truck



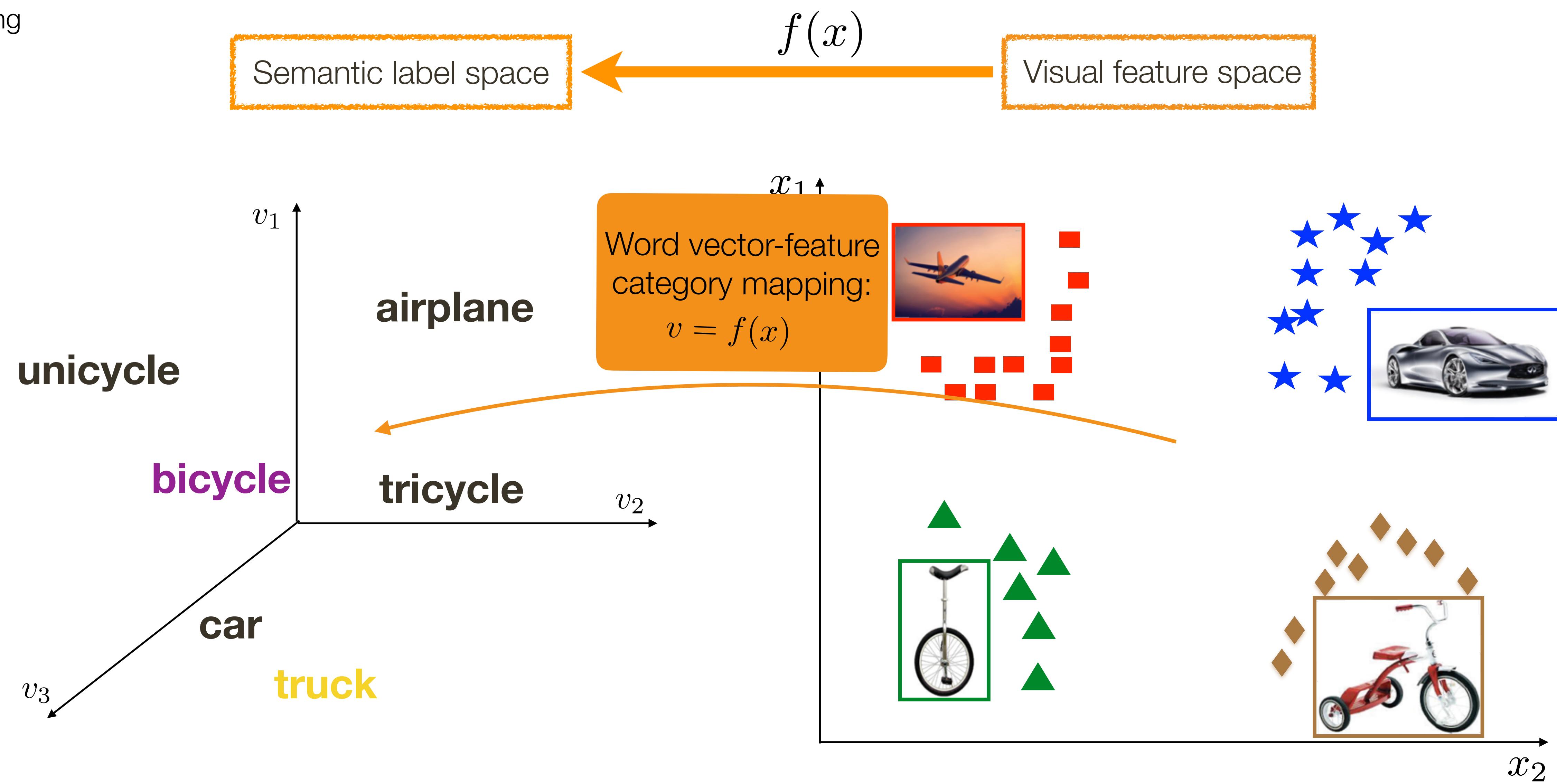
# Zero-shot Learning

Learning



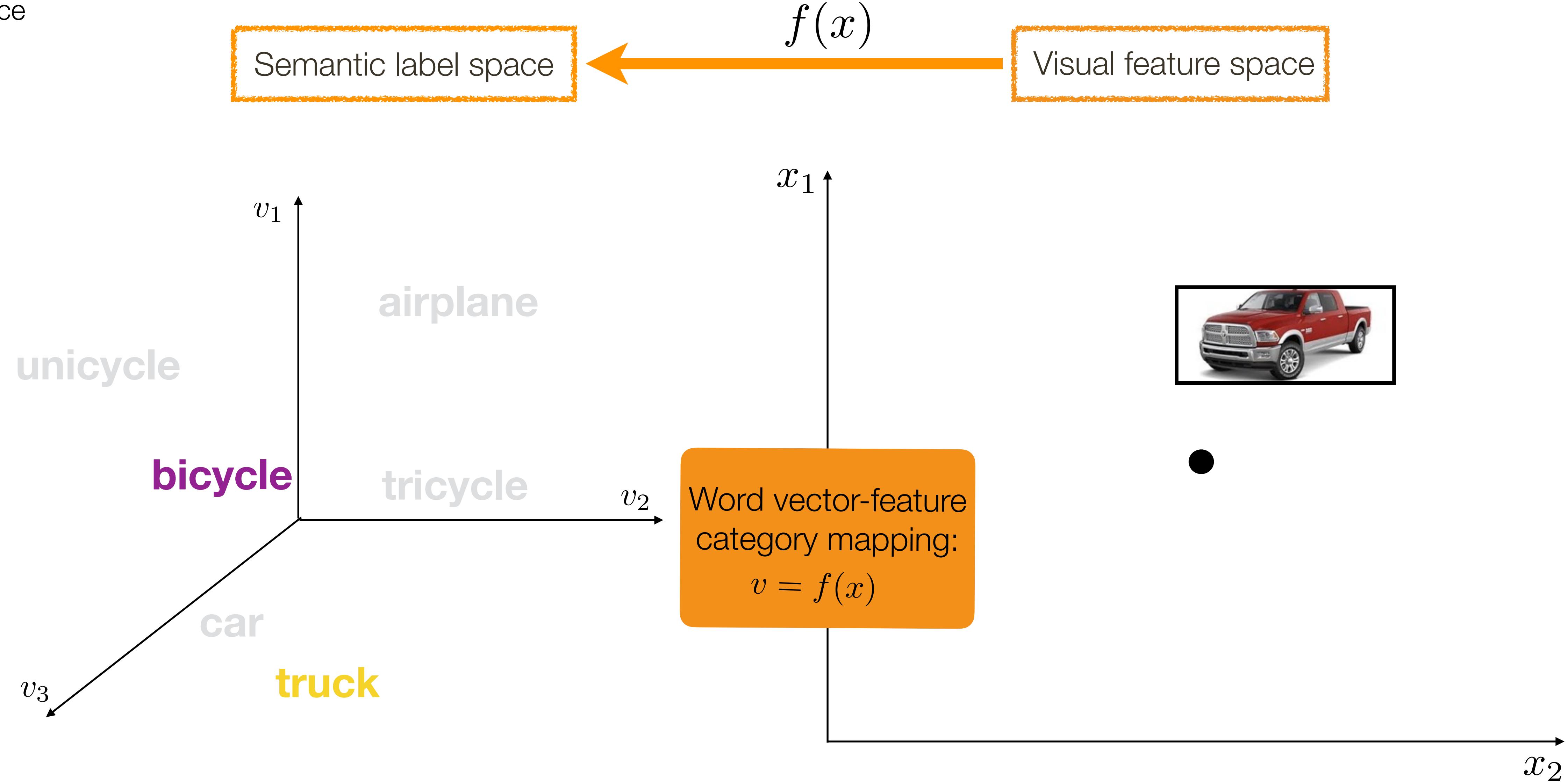
# Zero-shot Learning

Learning



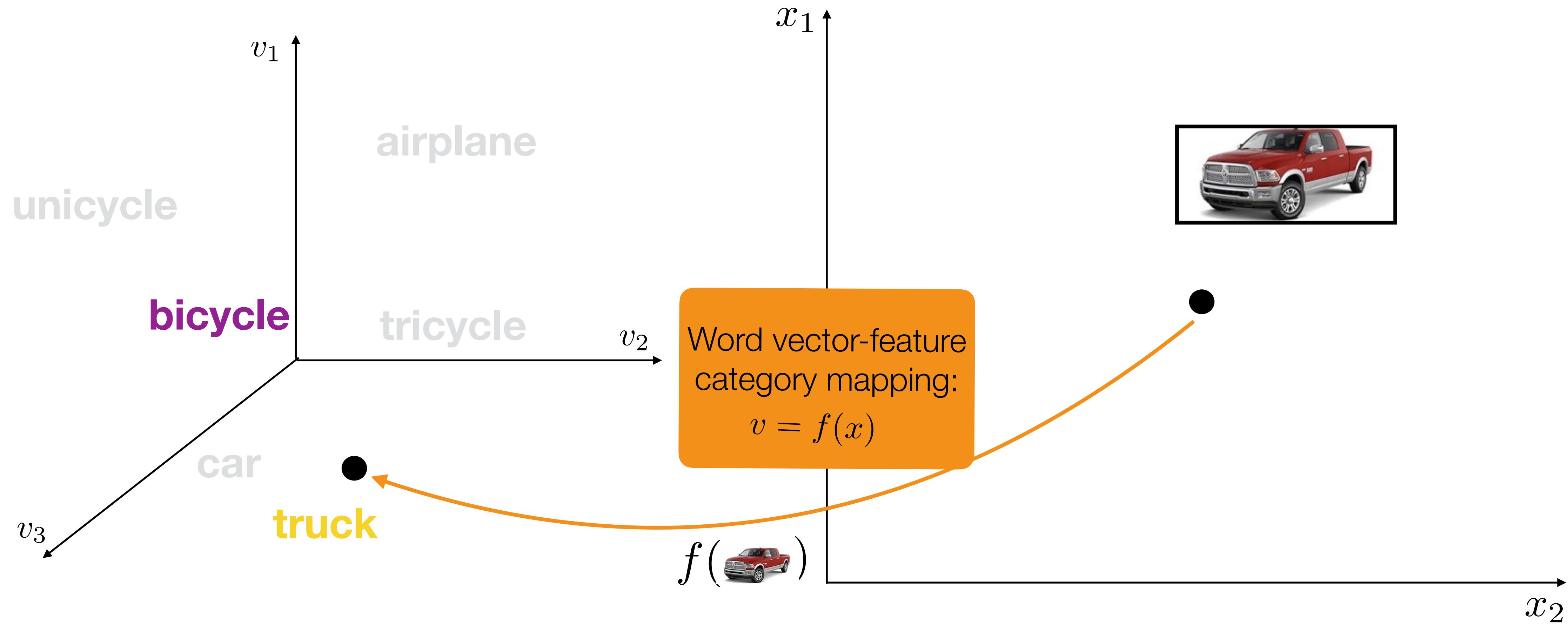
# Zero-shot Learning

Inference



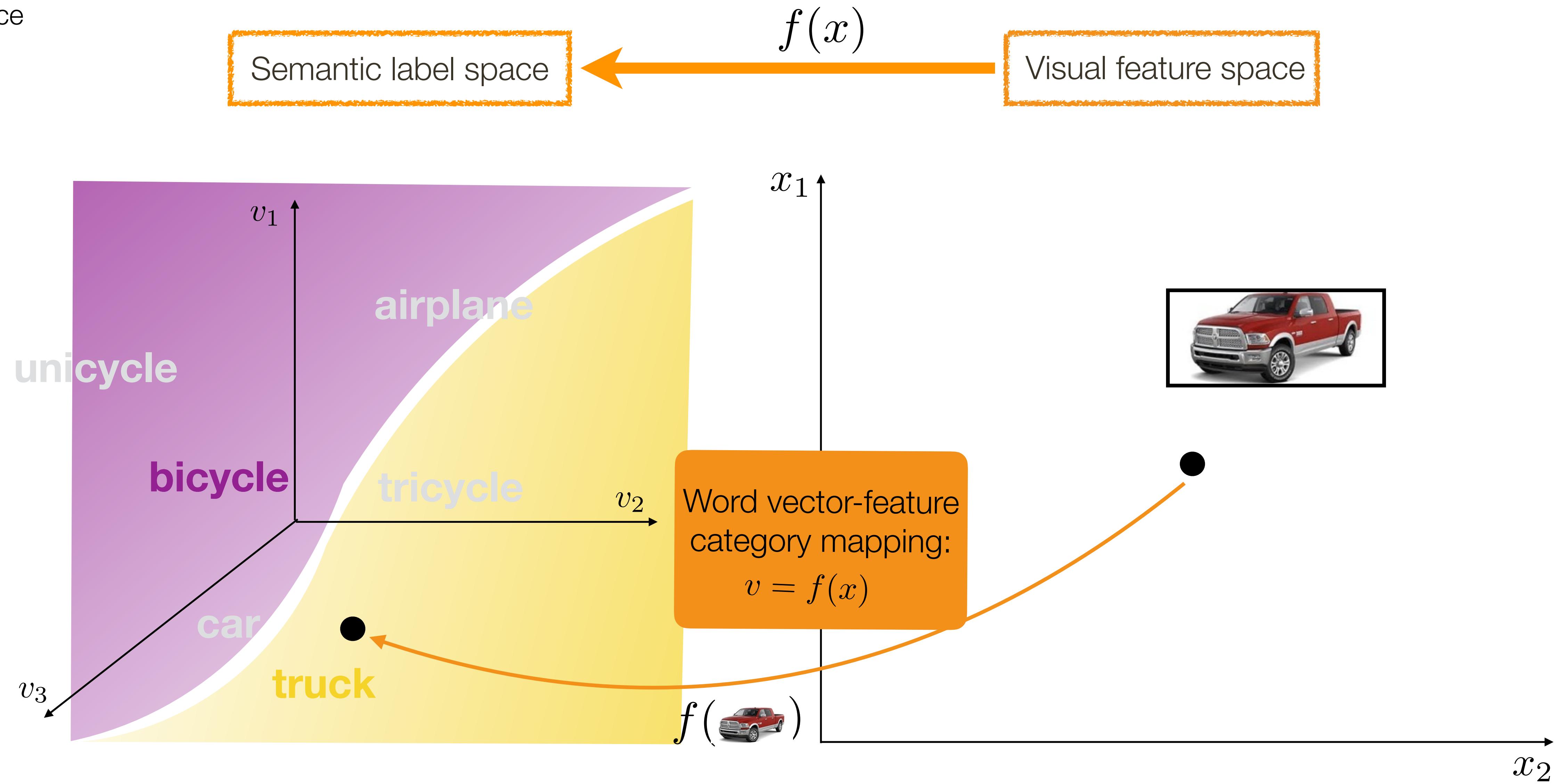
# Zero-shot Learning

Inference



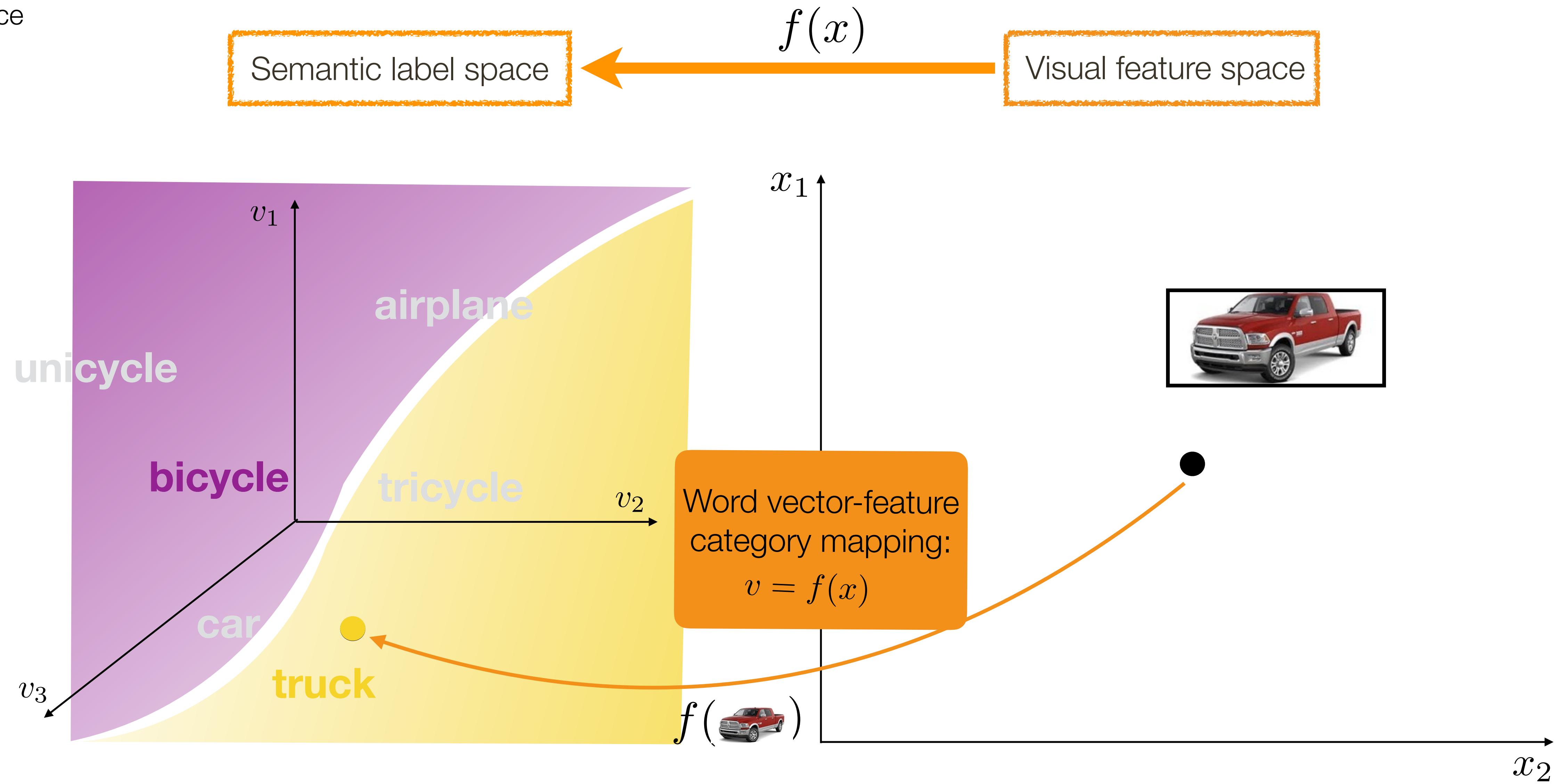
# Zero-shot Learning

Inference



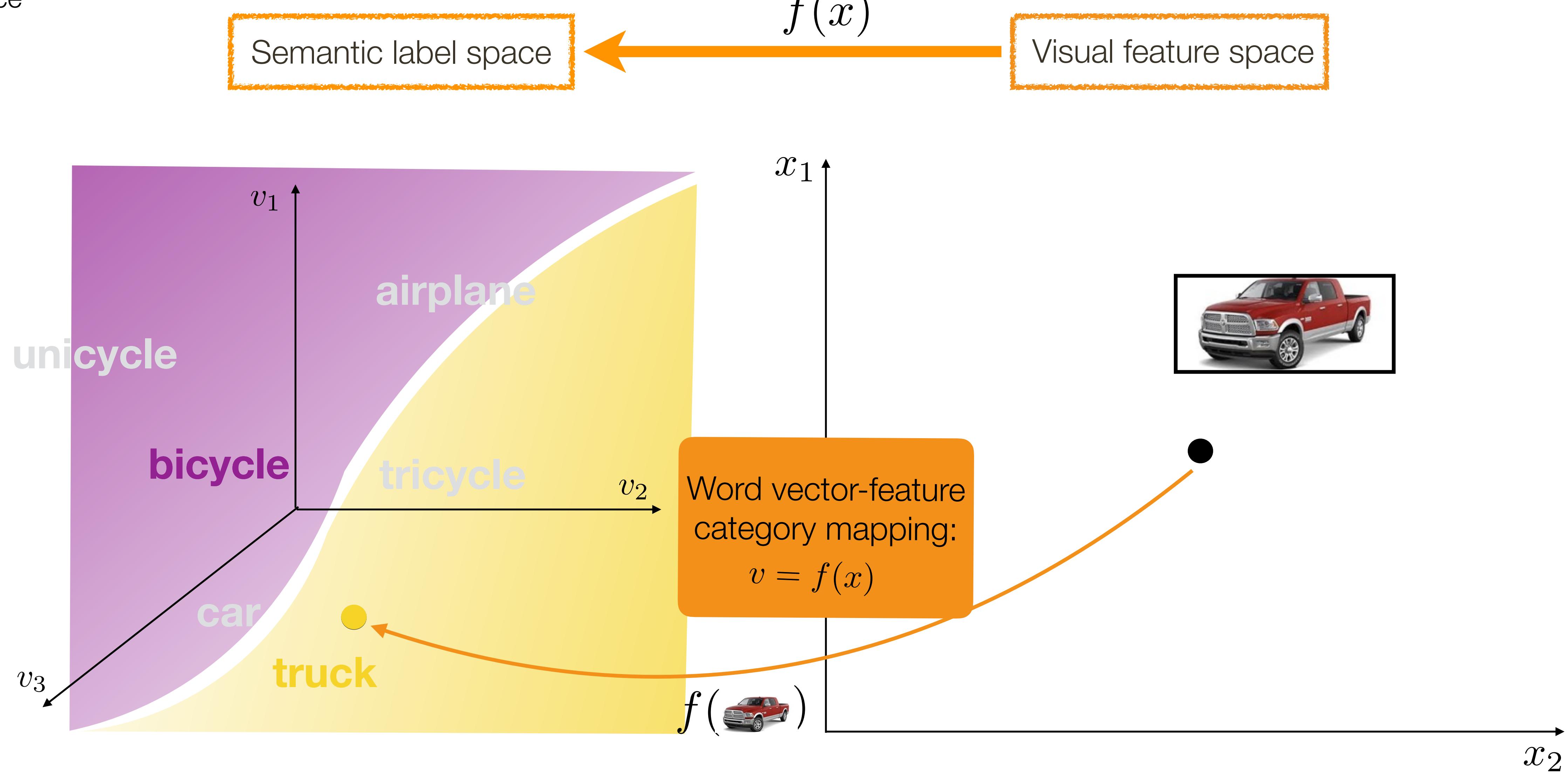
# Zero-shot Learning

Inference



# Zero-shot Learning

Inference



**Key Question:** How do we define semantic space?

# Semantic Label Vector Spaces

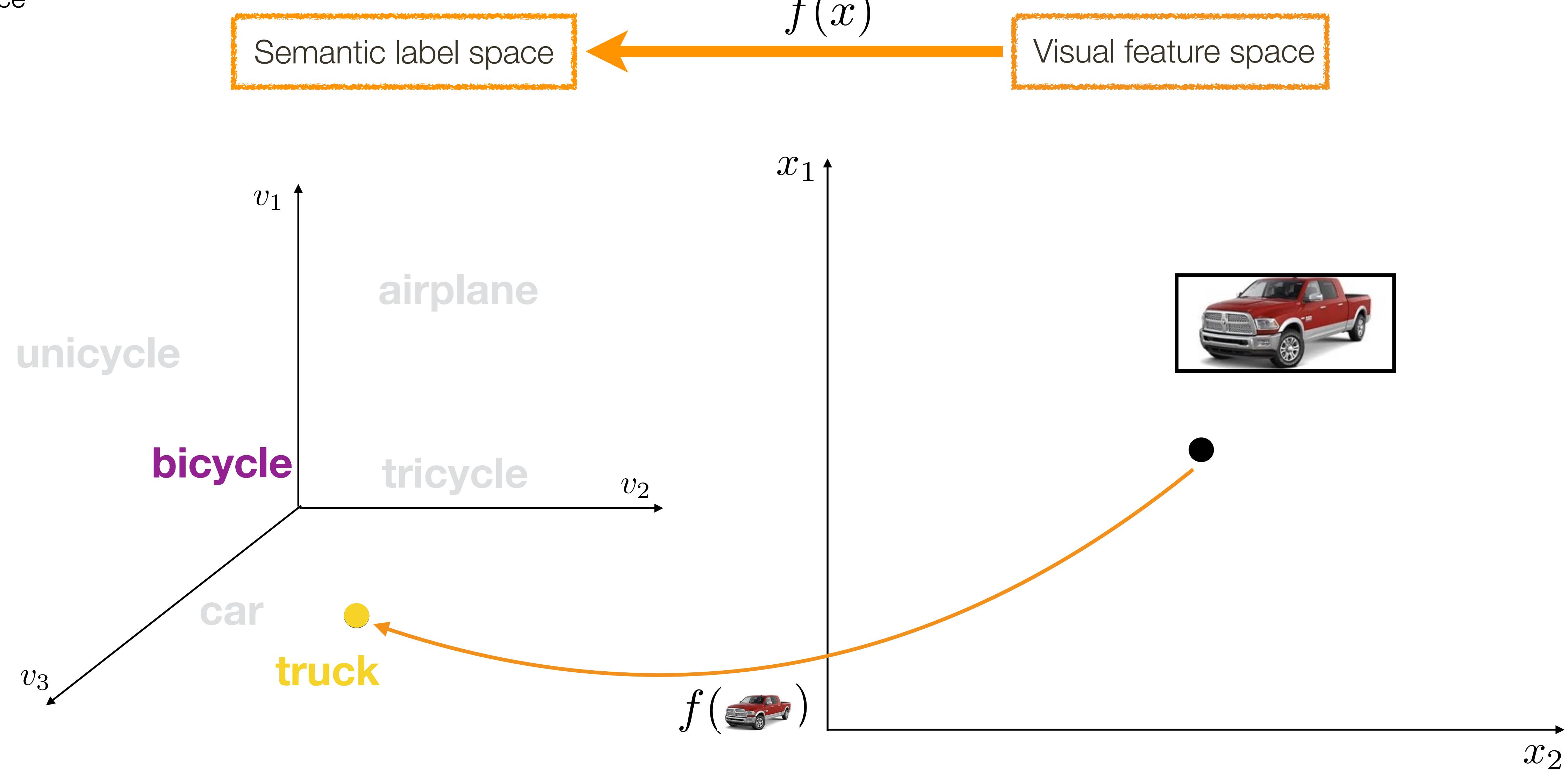
Spaces	Type	Advantages	Disadvantages
Semantic Attributes	Supervised	Good interpretability of each dimension: airplane := fixed_wing, propelled, has_pilot	Manual annotation Limited vocabulary
<b>Semantic Word Vectors</b> (e.g. word2vec)	Unsupervised	Good vector representation for millions of vocabulary $v(\text{Berlin}) - v(\text{Germany}) = v(\text{Paris}) - v(\text{France})$	Limited interpretability of each dimension

# Semantic Label Vector Spaces

Spaces	Type	Advantages	Disadvantages
Semantic Attributes	Supervised	Good interpretability of each dimension: airplane := fixed_wing, propelled, has_pilot	Manual annotation Limited vocabulary
Semantic Word Vectors (e.g. word2vec)	Unsupervised	Good vector representation for millions of vocabulary $v(\text{Berlin}) - v(\text{Germany}) = v(\text{Paris}) - v(\text{France})$	Limited interpretability of each dimension

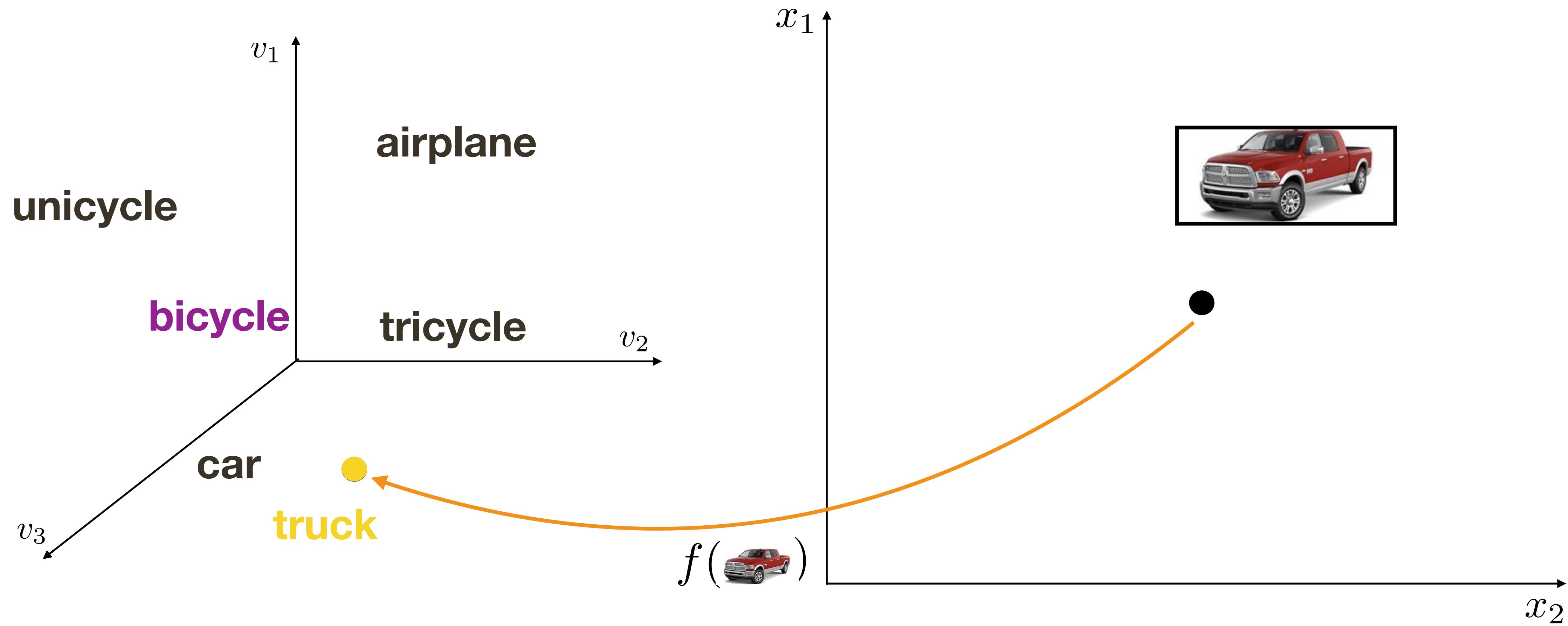
# Zero-shot Learning

Inference



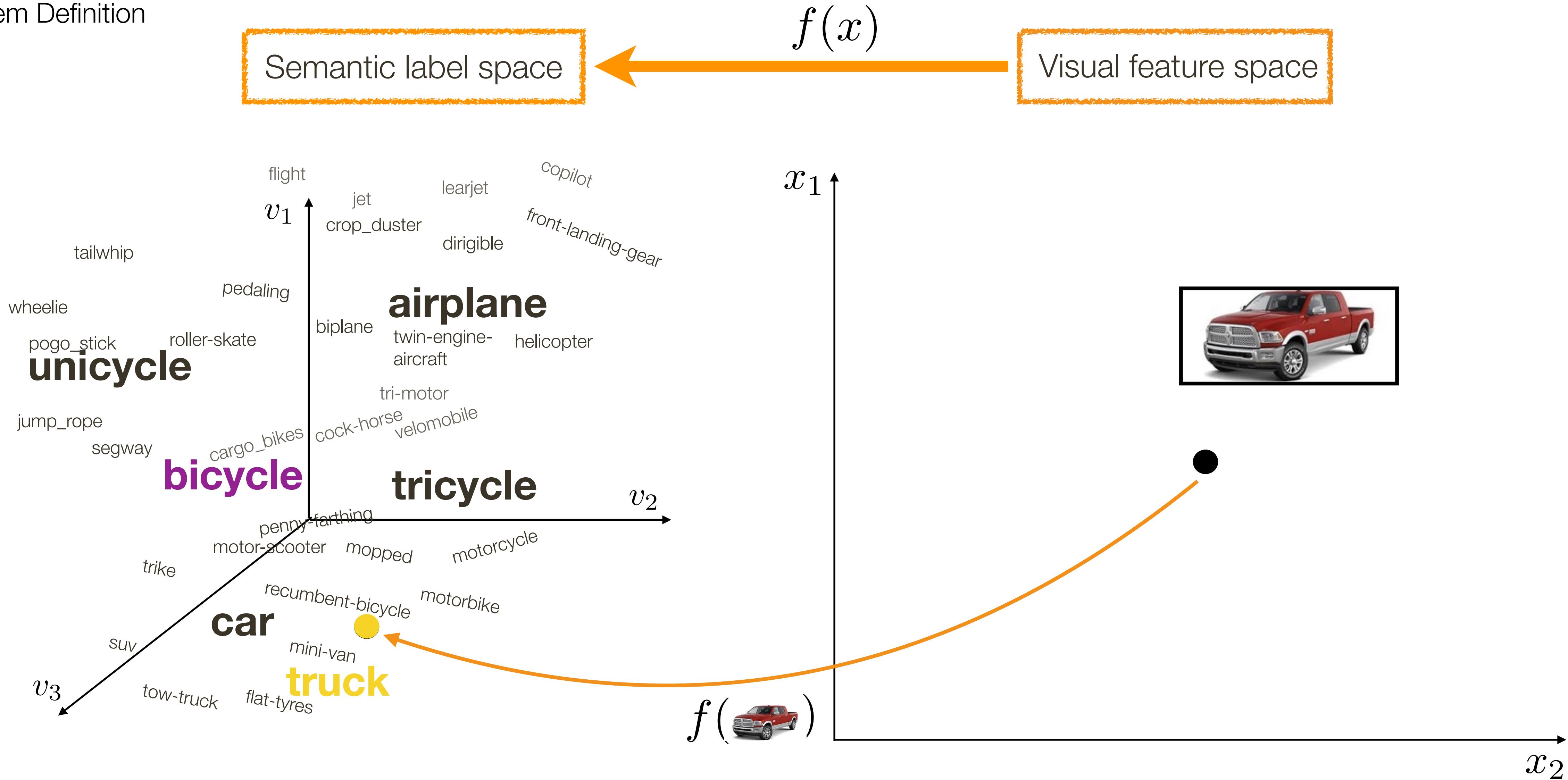
# Open-set Recognition

Problem Definition



# Open-set Recognition

Problem Definition



# Summary:

**Supervised Learning:** [Fei-Fei *et al.* TPAMI'06], [Deng *et al.* ECCV'14], [Torralba *et al.* TPAMI'08 ], [Weston *et al.* IJCAI'11]

Pros: Very good quantitative performance

Cons: Relatively small vocabulary (~1,000 classes)  
Requires ***manual labeling*** of all the data

# Summary:

## Supervised Learning: [Fei-Fei *et al.* TPAMI'06], [Deng *et al.* ECCV'14], [Torralba *et al.* TPAMI'08 ], [Weston *et al.* IJCAI'11]

Pros: Very good quantitative performance

Cons: Relatively small vocabulary (~1,000 classes)  
Requires ***manual labeling*** of all the data

## Zero-shot Learning: [Palatucci *et al.* NIPS'09],[Lampert *et al.* CVPR'09], [Farhadi *et al.* CVPR'09], [Rohrbach *et al.* CVPR'10]

Pros: Does not require instance labeling for target classes

Cons: Typically limited to recognition with target classes only  
Relatively ***small vocabulary*** (~50-200 classes typically)

# Summary:

## Supervised Learning: [Fei-Fei *et al.* TPAMI'06], [Deng *et al.* ECCV'14], [Torralba *et al.* TPAMI'08 ], [Weston *et al.* IJCAI'11]

Pros: Very good quantitative performance

Cons: Relatively small vocabulary (~1,000 classes)  
Requires ***manual labeling*** of all the data

## Zero-shot Learning: [Palatucci *et al.* NIPS'09],[Lampert *et al.* CVPR'09], [Farhadi *et al.* CVPR'09], [Rohrbach *et al.* CVPR'10]

Pros: Does not require instance labeling for target classes

Cons: Typically limited to recognition with target classes only  
Relatively ***small vocabulary*** (~50-200 classes typically)

## Open-set Learning: [Scheirer *et al.* TPAMI'13], [Sattar *et al.* CVPR'15], [Bendale *et al.* CVPR'15] [Guadarrama *et al.* RSS'14]

Pros: Does not require instance labeling for target classes  
Large vocabulary (up to 310K classes)

# Summary:

## Supervised Learning: [Fei-Fei *et al.* TPAMI'06], [Deng *et al.* ECCV'14], [Torralba *et al.* TPAMI'08 ], [Weston *et al.* IJCAI'11]

Pros: Very good quantitative performance

Cons: Relatively small vocabulary (~1,000 classes)  
Requires **manual labeling** of all the data

## Zero-shot Learning: [Palatucci *et al.* NIPS'09],[Lampert *et al.* CVPR'09], [Farhadi *et al.* CVPR'09], [Rohrbach *et al.* CVPR'10]

Pros: Does not require instance labeling for target classes

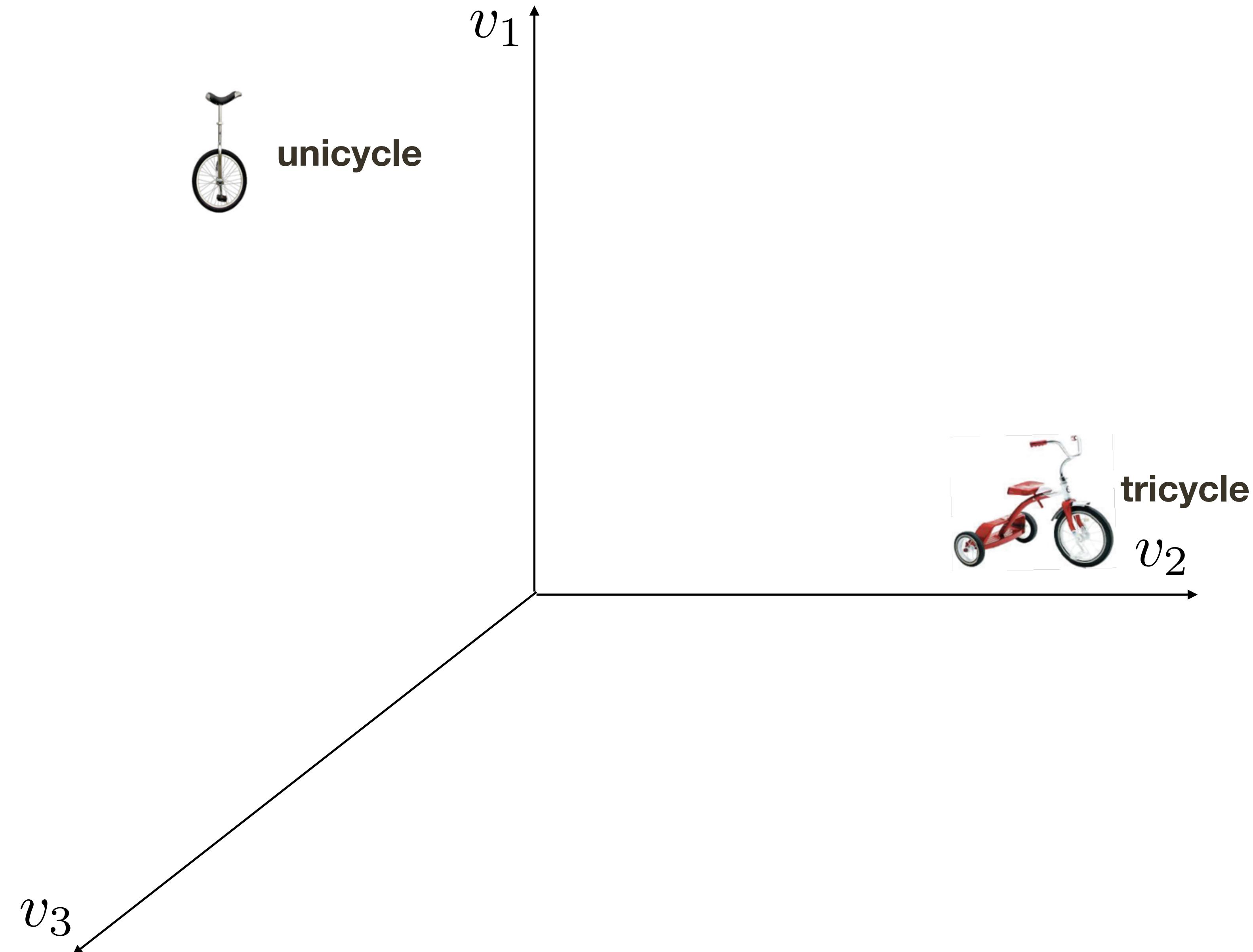
Cons: Typically limited to recognition with target classes only  
Relatively **small vocabulary** (~50-200 classes typically)

## Open-set Learning: [Scheirer *et al.* TPAMI'13], [Sattar *et al.* CVPR'15], [Bendale *et al.* CVPR'15] [Guadarrama *et al.* RSS'14]

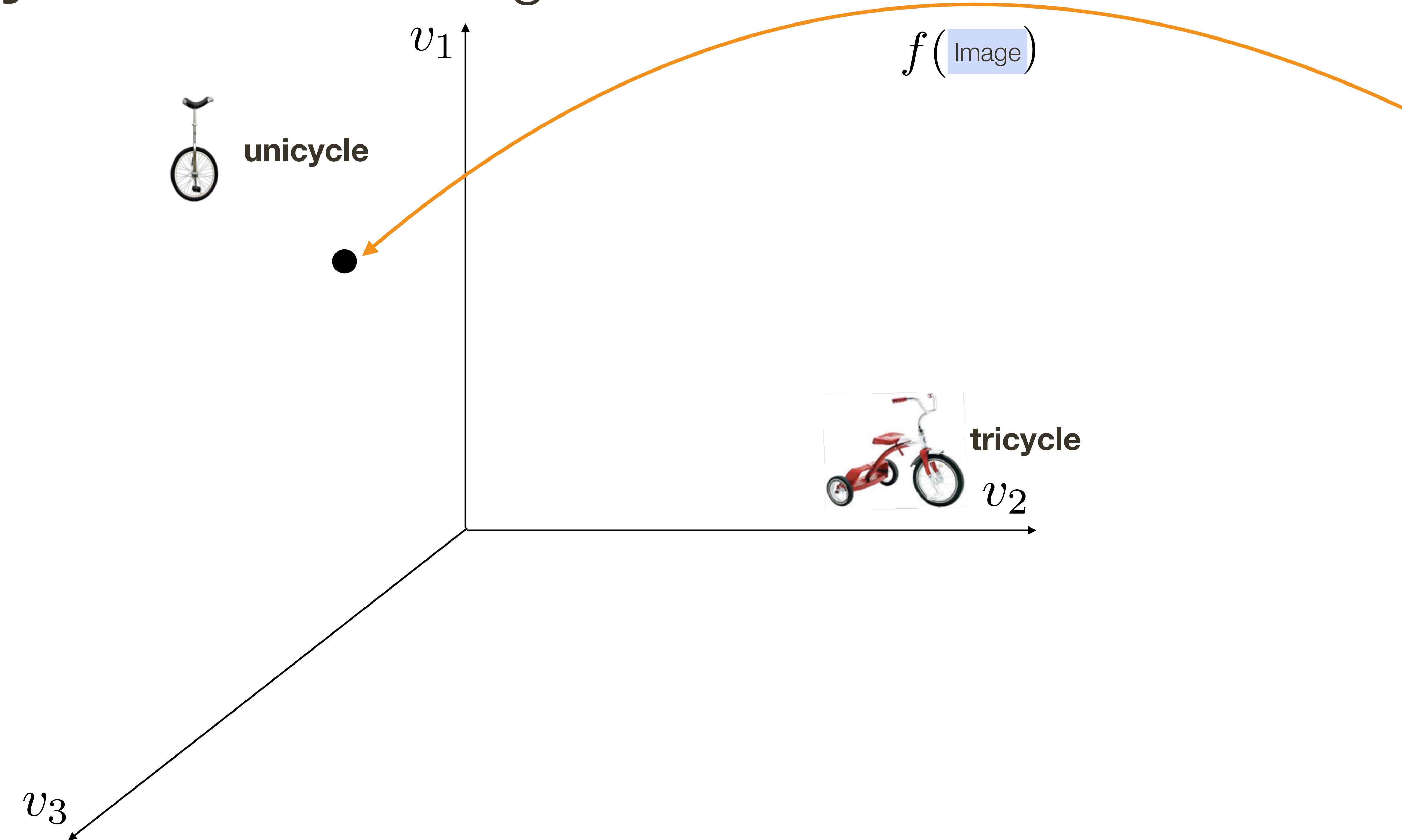
Pros: Does not require instance labeling for target classes  
Large vocabulary (up to 310K classes)

**Key Question:** Can this large vocabulary be actually useful for recognition?

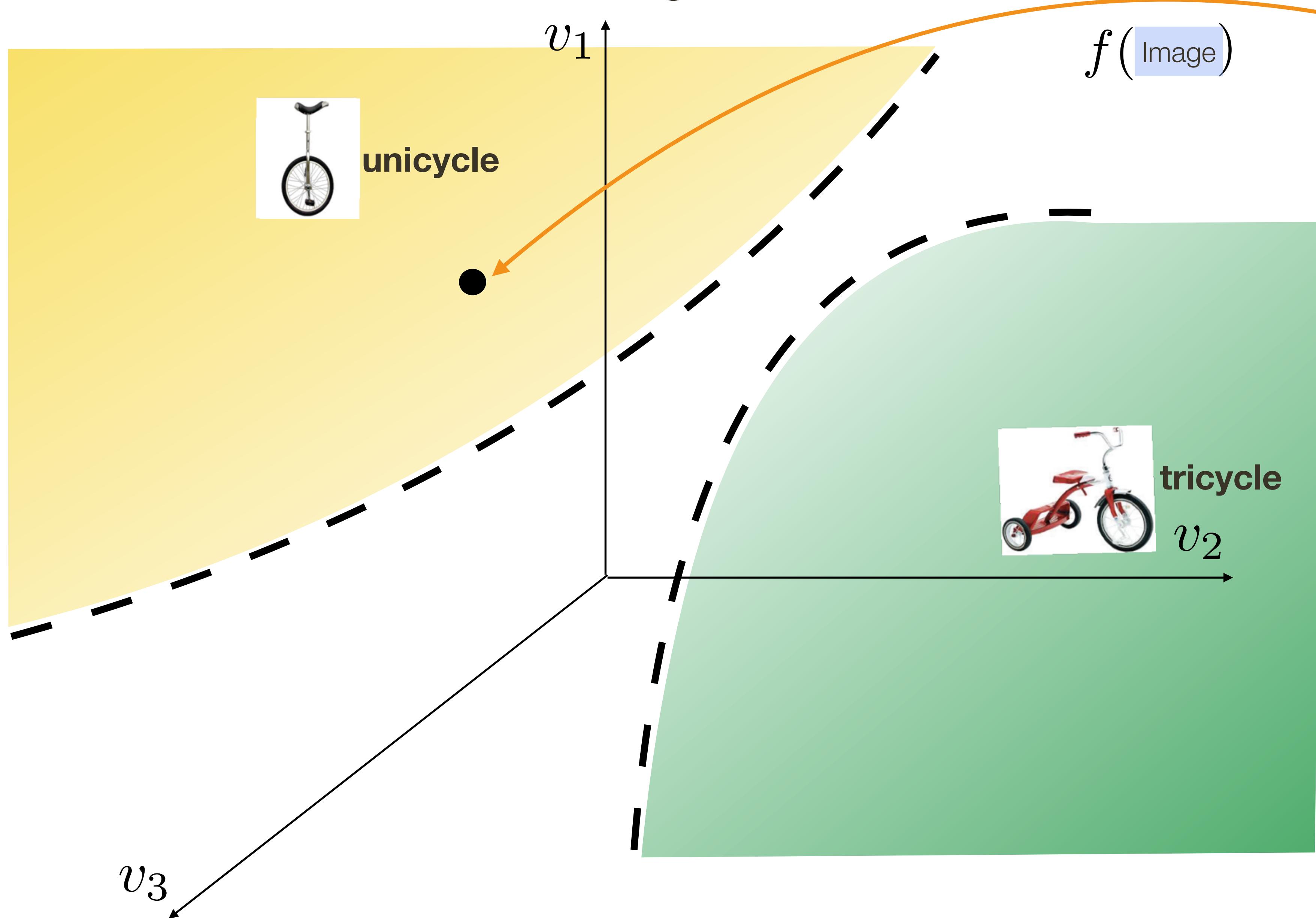
# Vocabulary-Informed Recognition



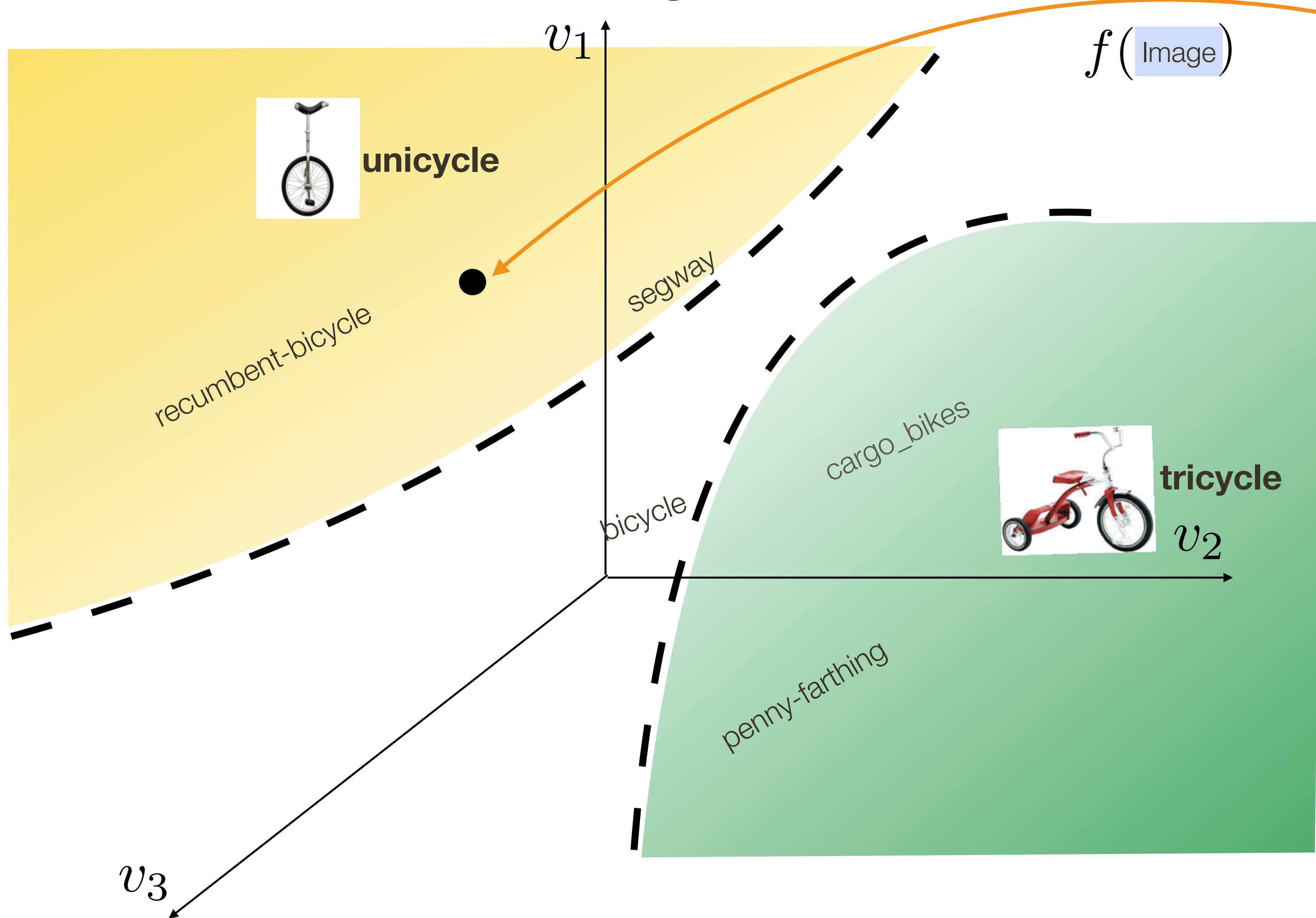
# Vocabulary-Informed Recognition



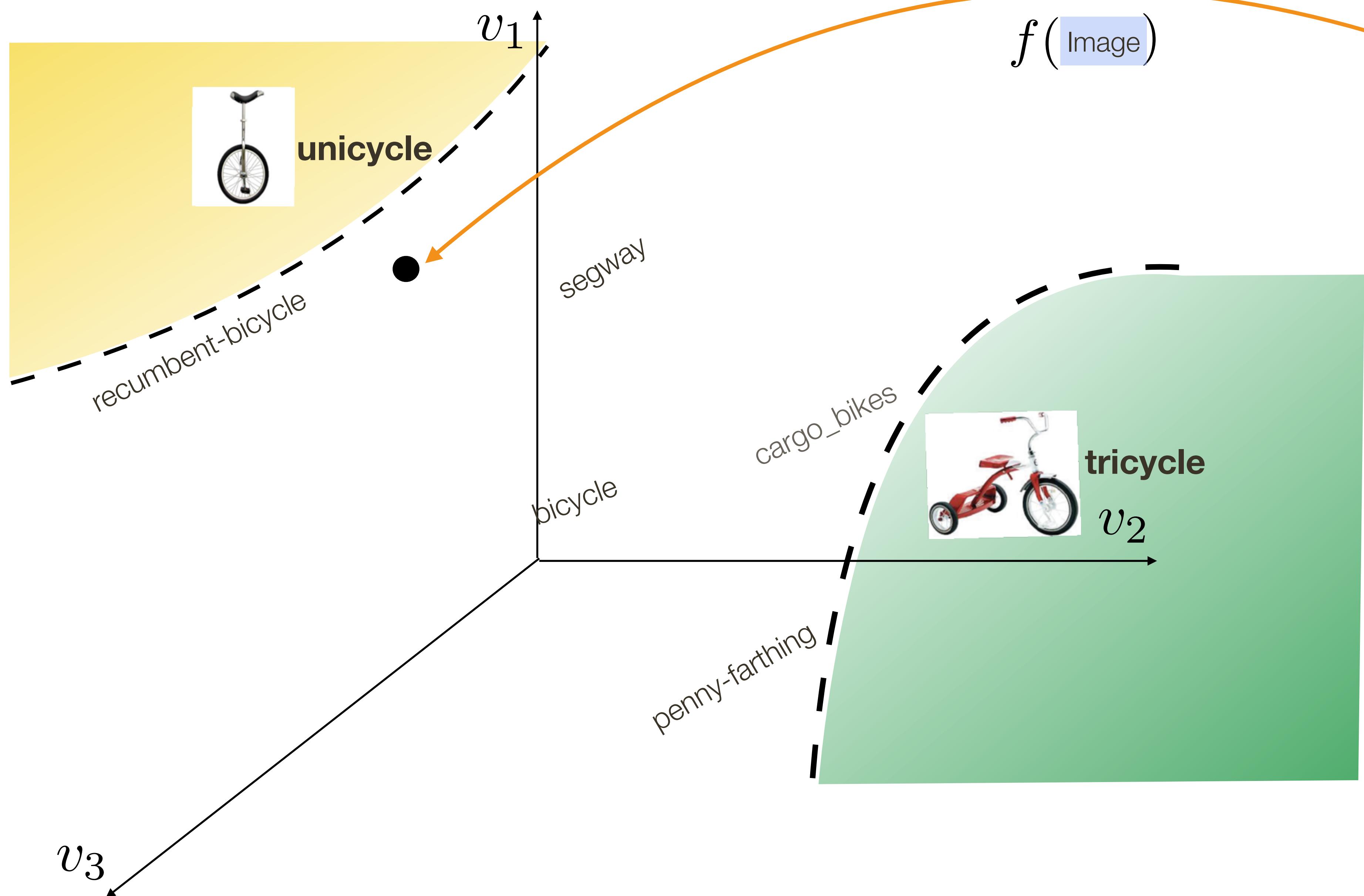
# Vocabulary-Informed Recognition



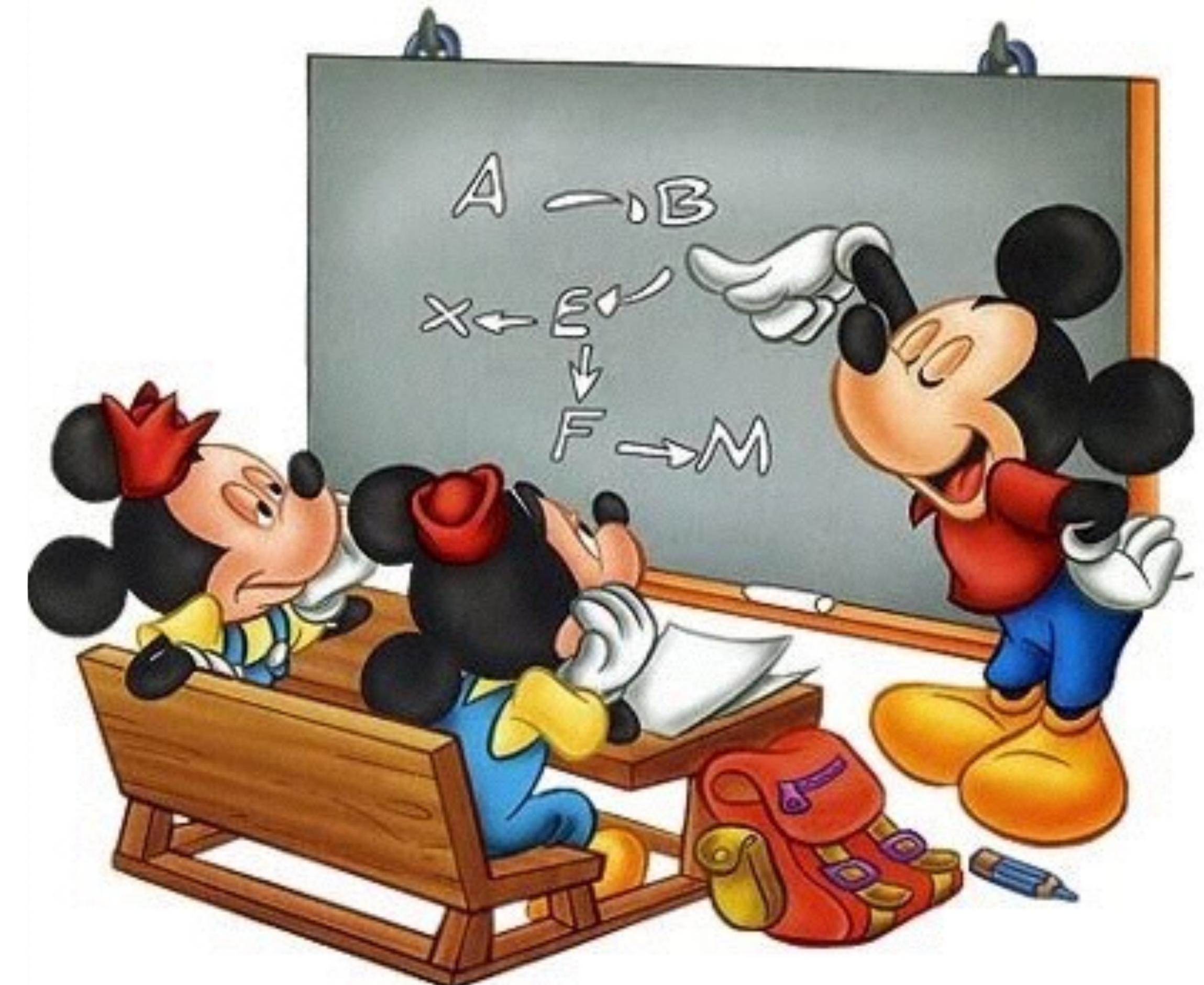
# Vocabulary-Informed Recognition



# Vocabulary-Informed Recognition

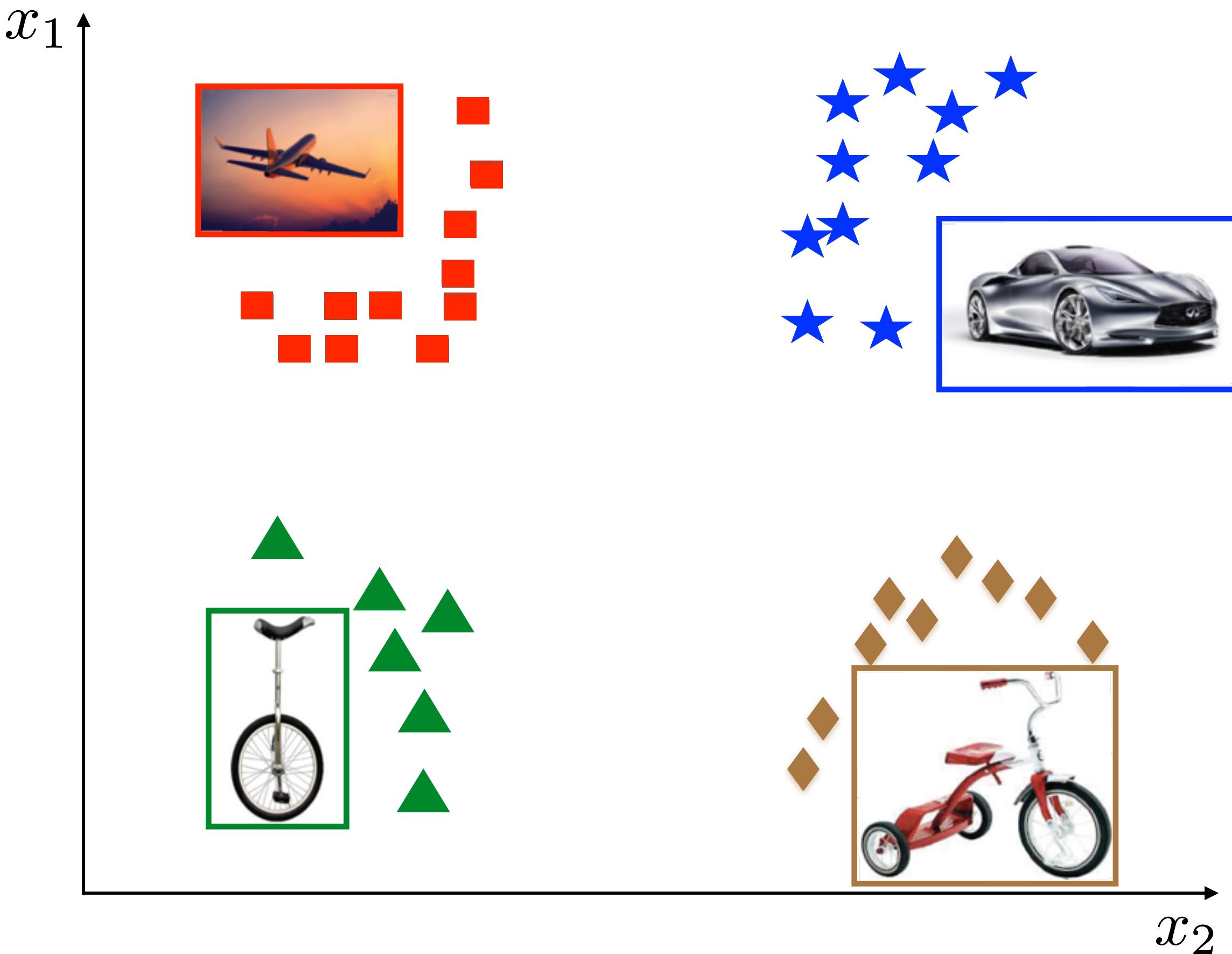
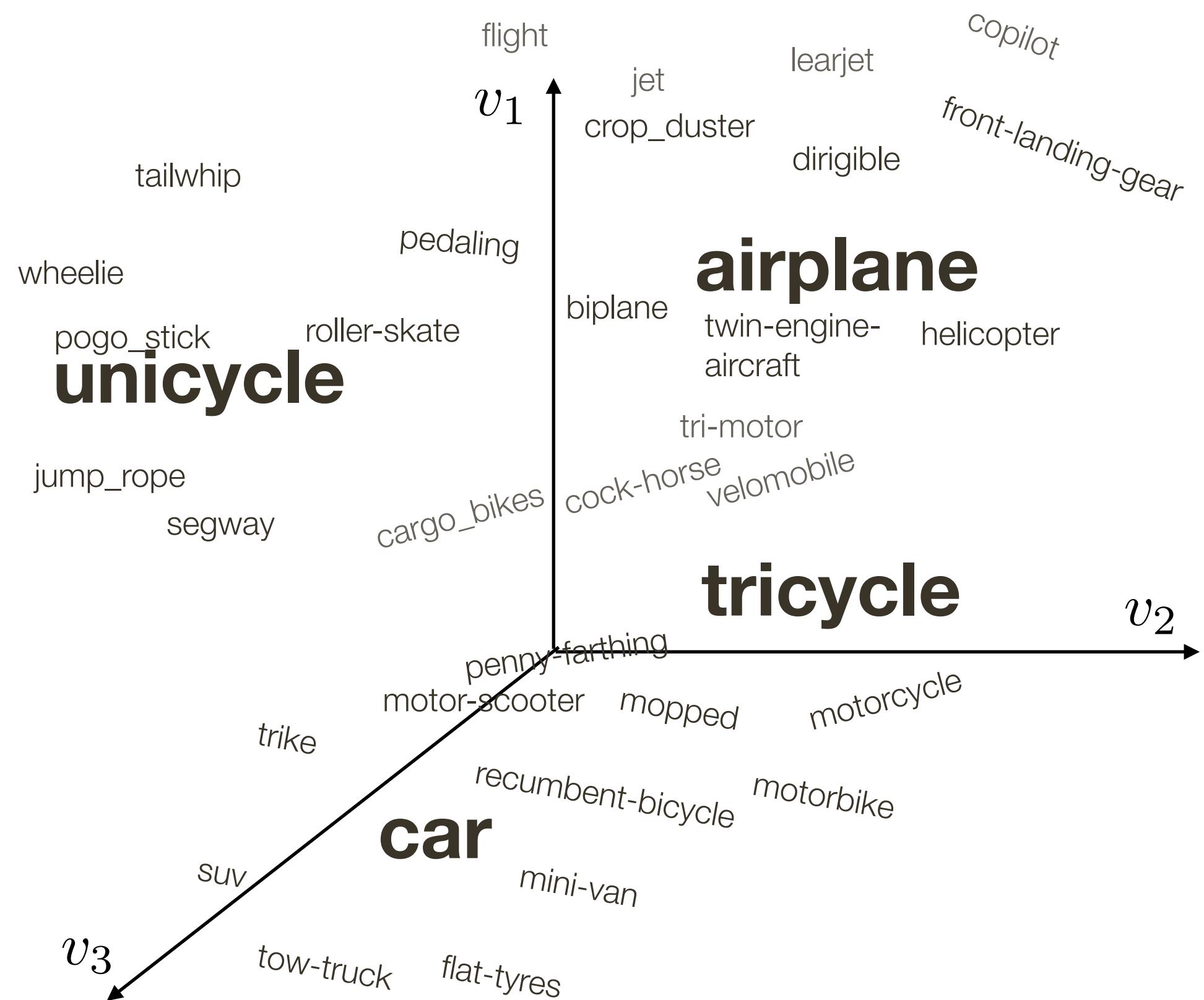
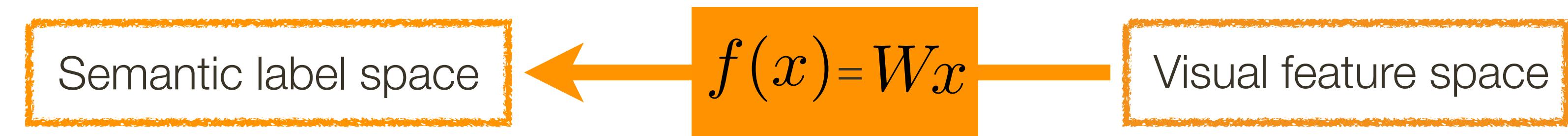


# Formulation



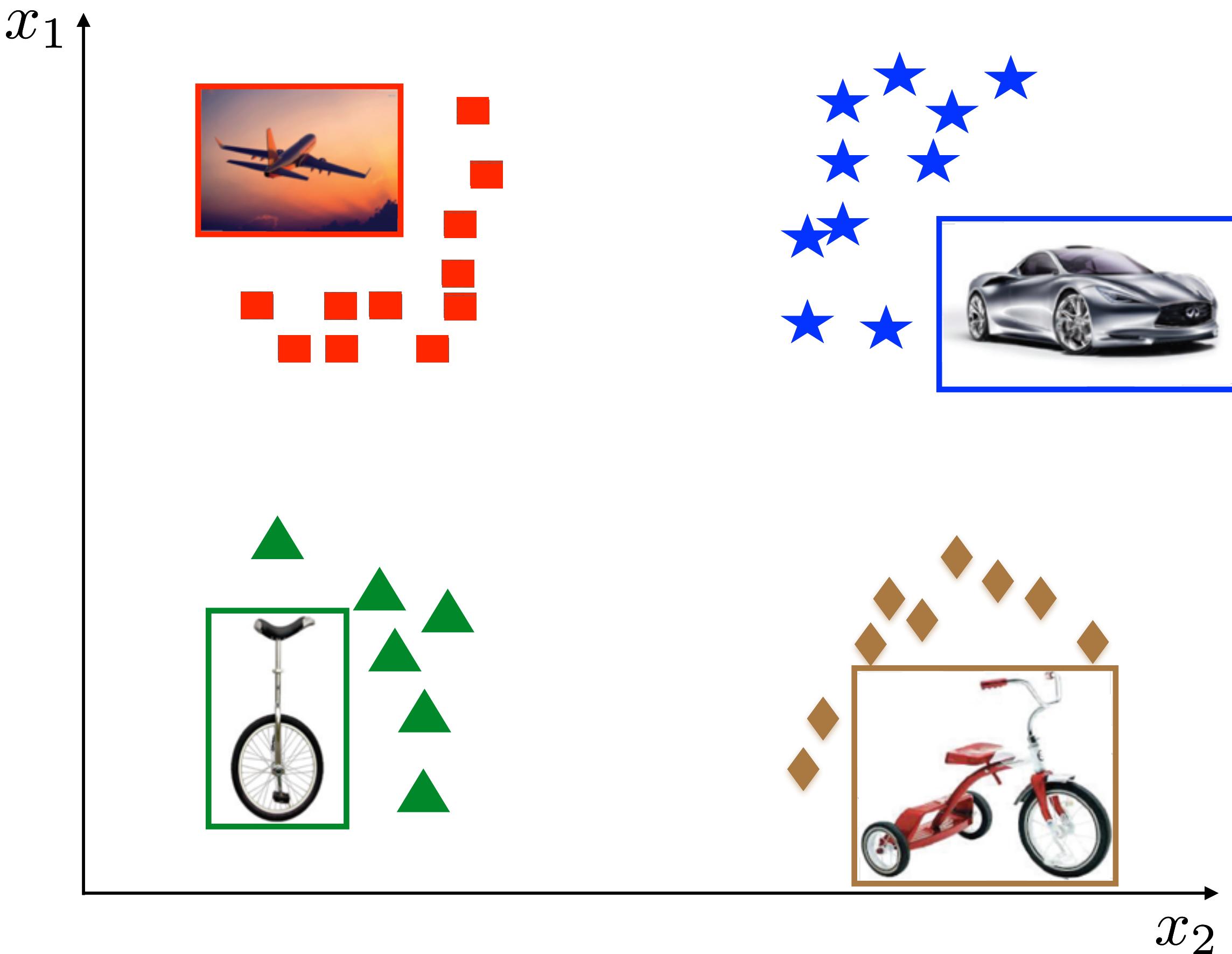
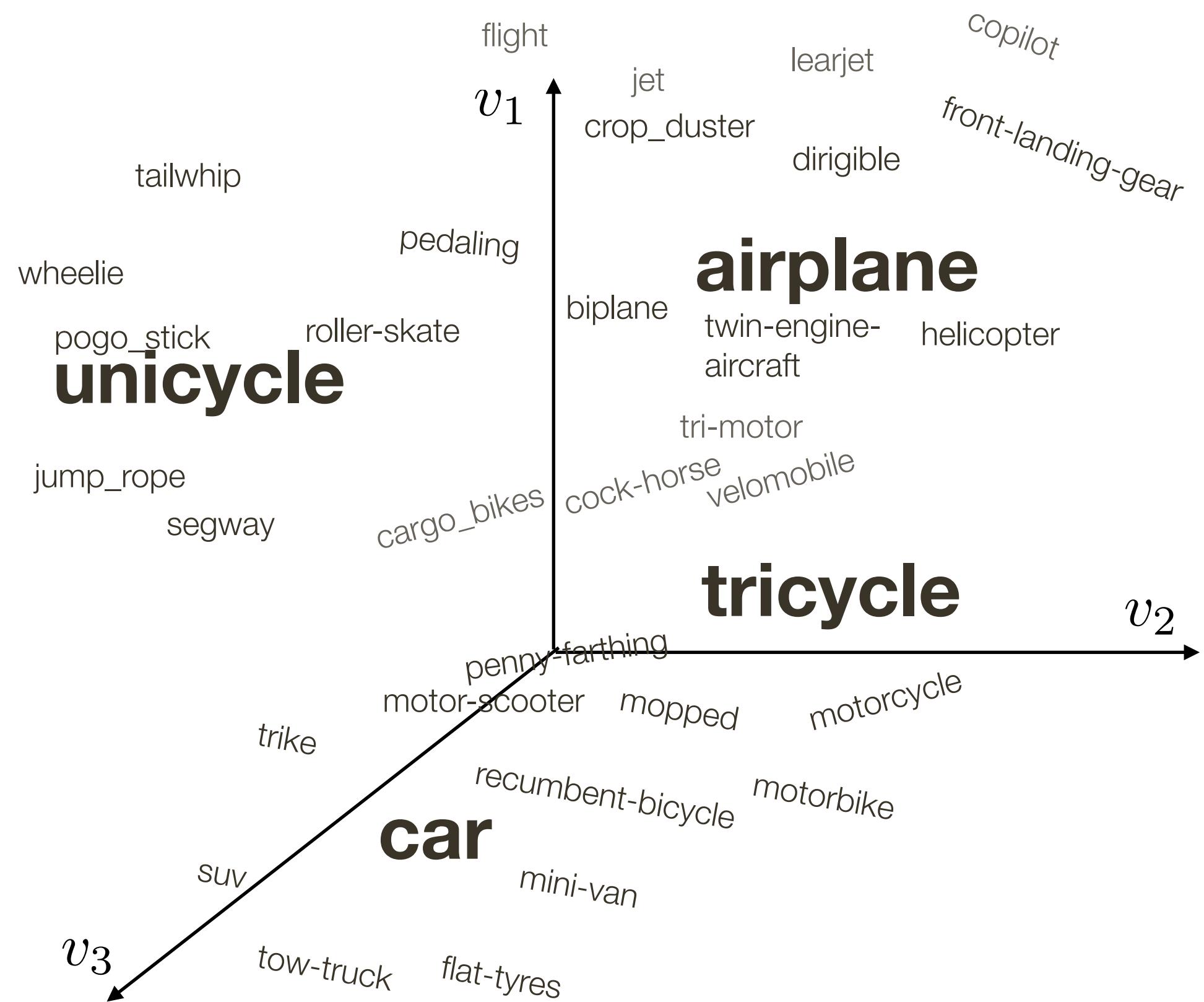
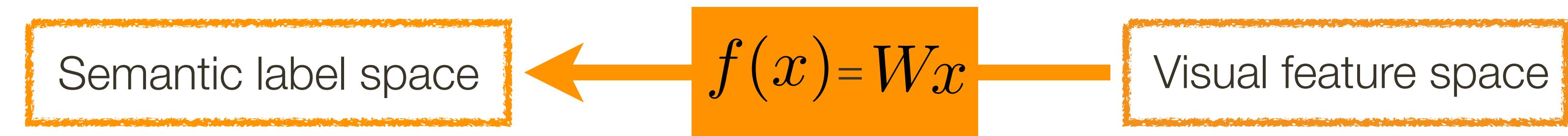
# Semi-Supervised Vocabulary-Informed Learning (SS-Voc)

Regression term



# Semi-Supervised Vocabulary-Informed Learning (SS-Voc)

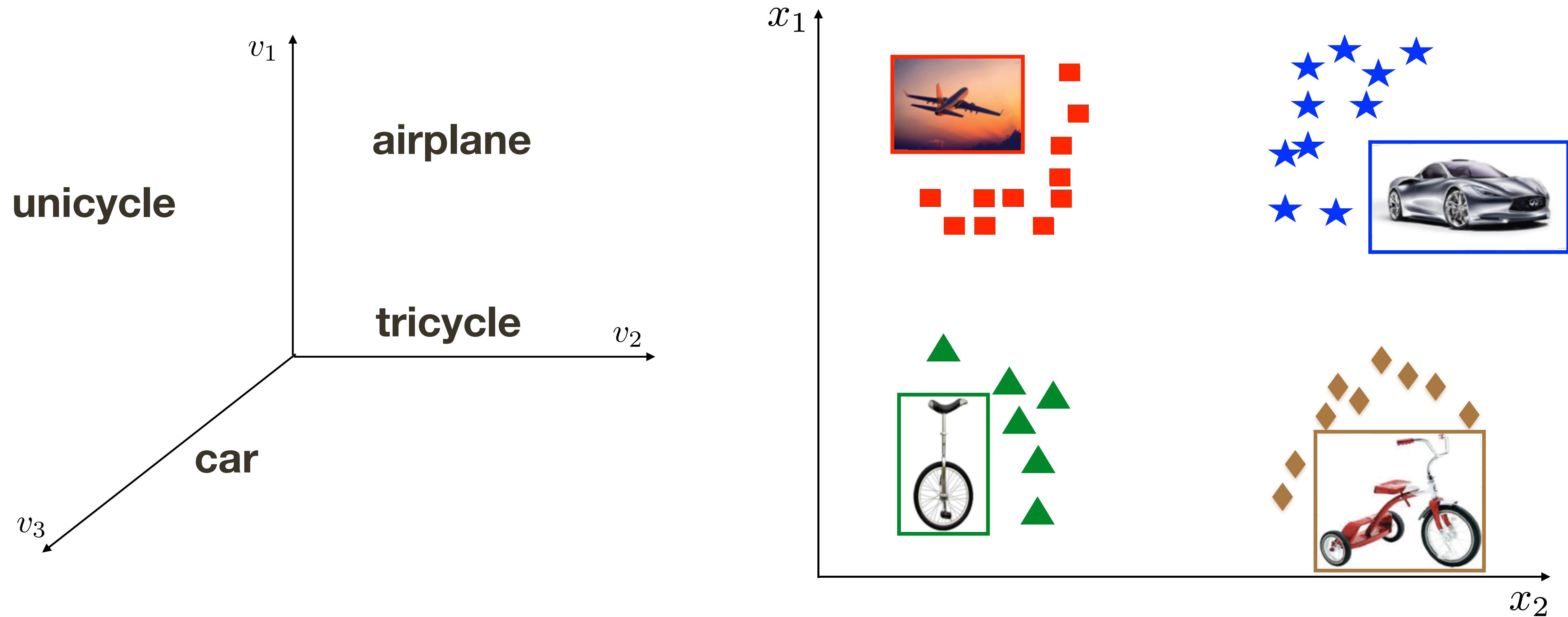
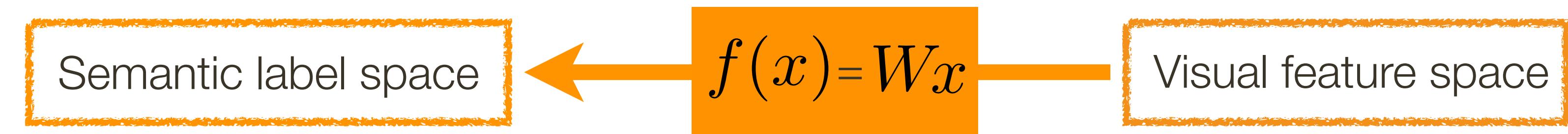
Regression term



$$\sum \mathcal{L}_\epsilon(f(\text{airplane image}, W), \mathbf{u}_{\text{airplane}})$$

# Semi-Supervised Vocabulary-Informed Learning (SS-Voc)

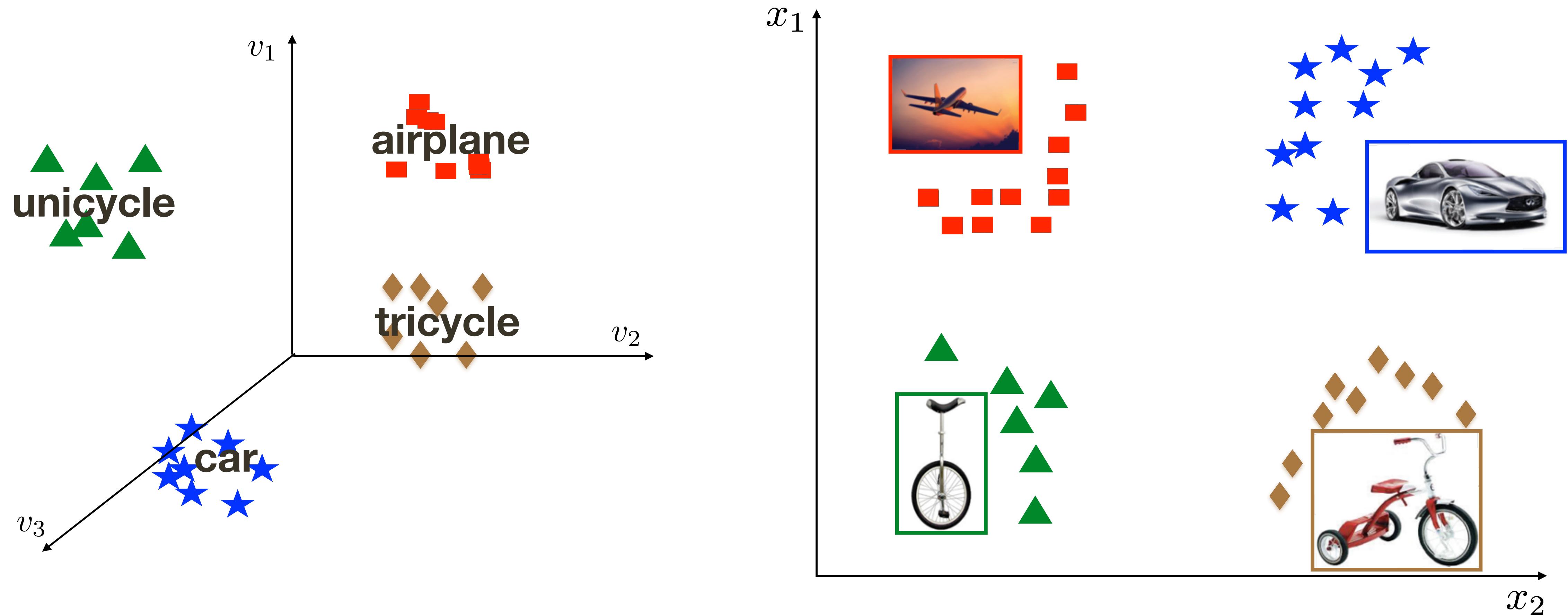
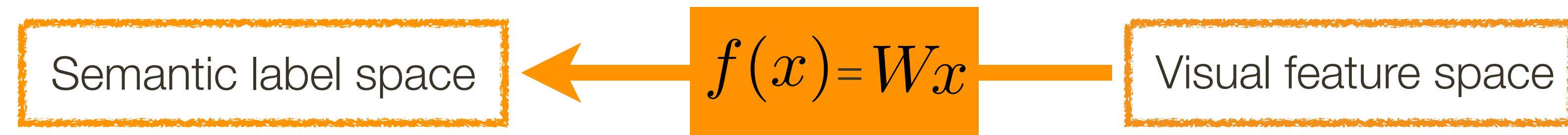
Regression term



$$\sum \mathcal{L}_\epsilon(f(\text{airplane image}, W), \mathbf{u}_{\text{airplane}})$$

# Semi-Supervised Vocabulary-Informed Learning (SS-Voc)

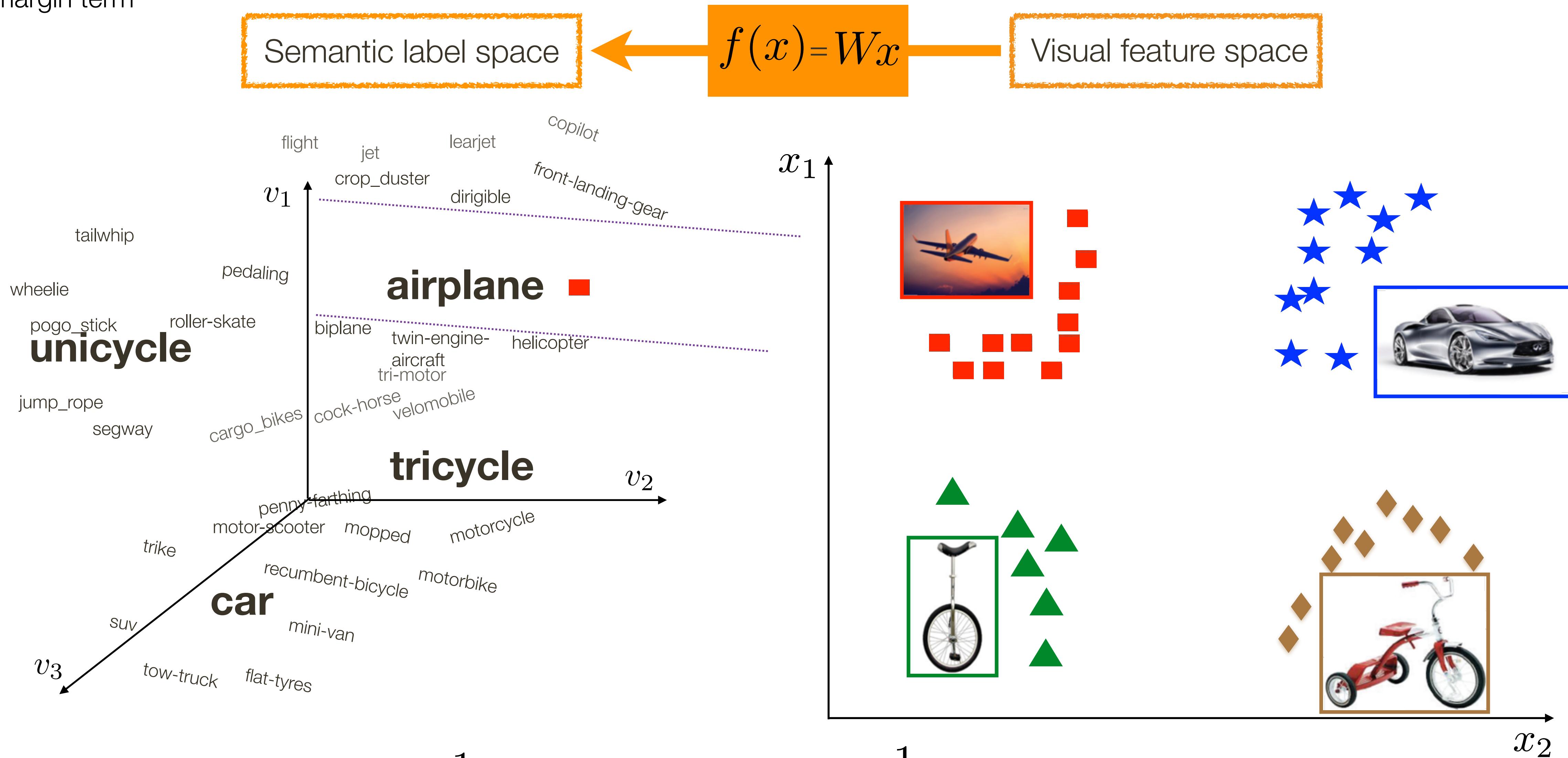
Regression term



$$\sum \mathcal{L}_\epsilon(f(\text{airplane}), \mathbf{u}_{\text{airplane}})$$

# Semi-Supervised Vocabulary-Informed Learning (SS-Voc)

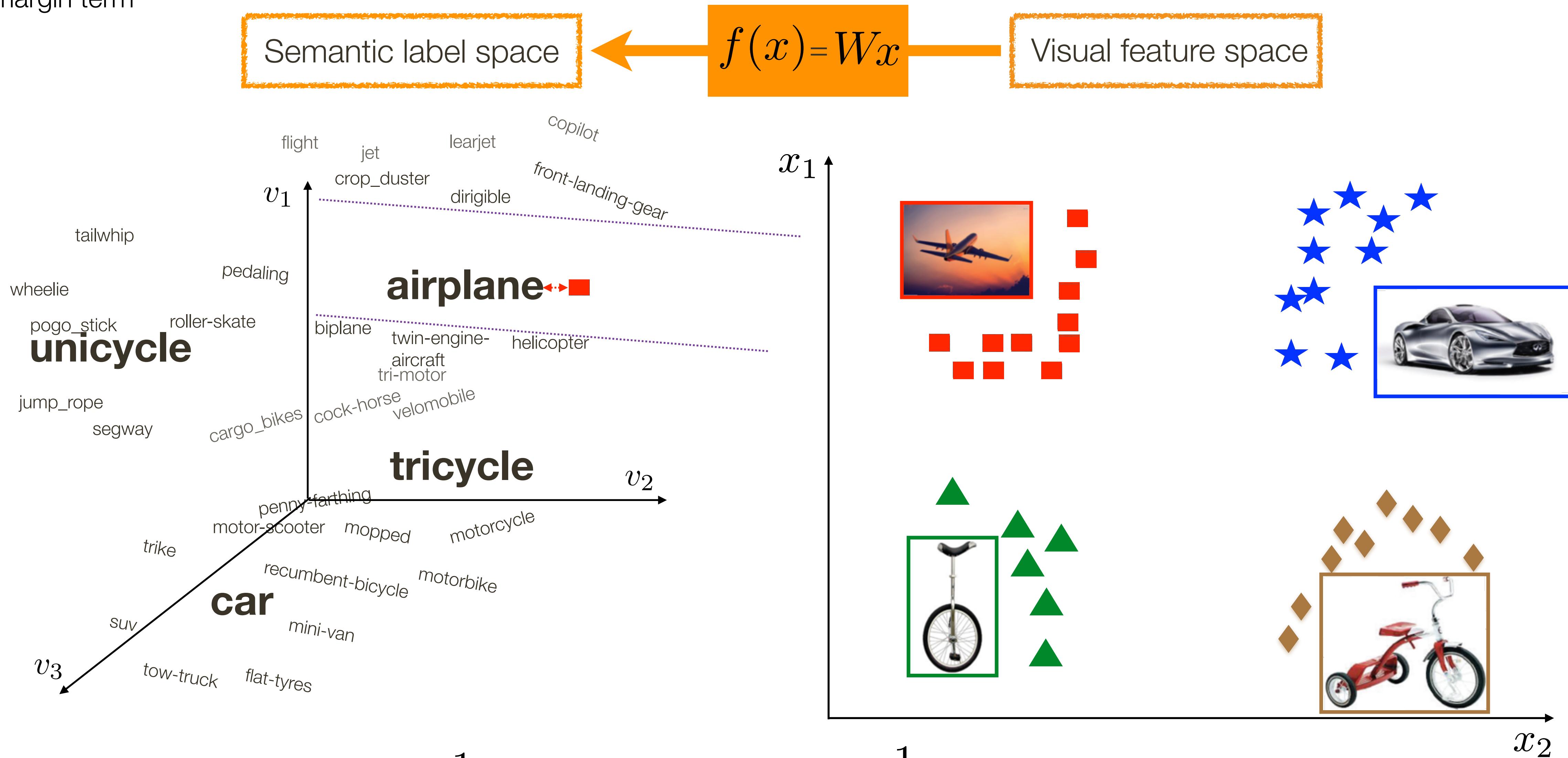
Max-margin term



$$\sum [C - \frac{1}{2} \mathcal{L}_\epsilon(f(\text{airplane}), \mathbf{u}_{\text{airplane}}) + \frac{1}{2} \mathcal{L}_\epsilon(f(\text{helicopter}), \mathbf{u}_{\text{helicopter}})]$$

# Semi-Supervised Vocabulary-Informed Learning (SS-Voc)

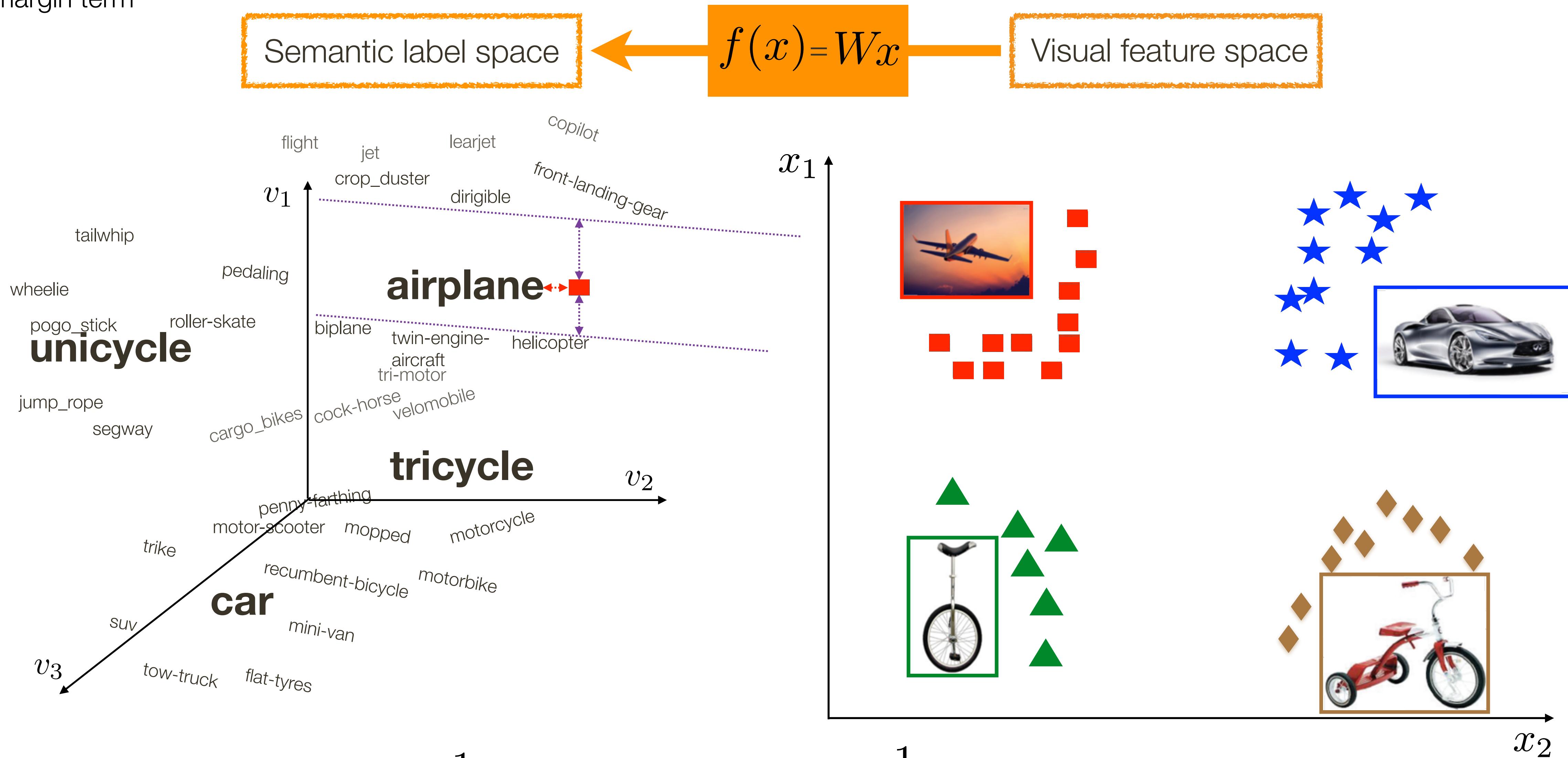
Max-margin term



$$\sum [C - \frac{1}{2} \mathcal{L}_\epsilon(f(\text{airplane}), \mathbf{u}_{\text{airplane}}) + \frac{1}{2} \mathcal{L}_\epsilon(f(\text{helicopter}), \mathbf{u}_{\text{helicopter}})]$$

# Semi-Supervised Vocabulary-Informed Learning (SS-Voc)

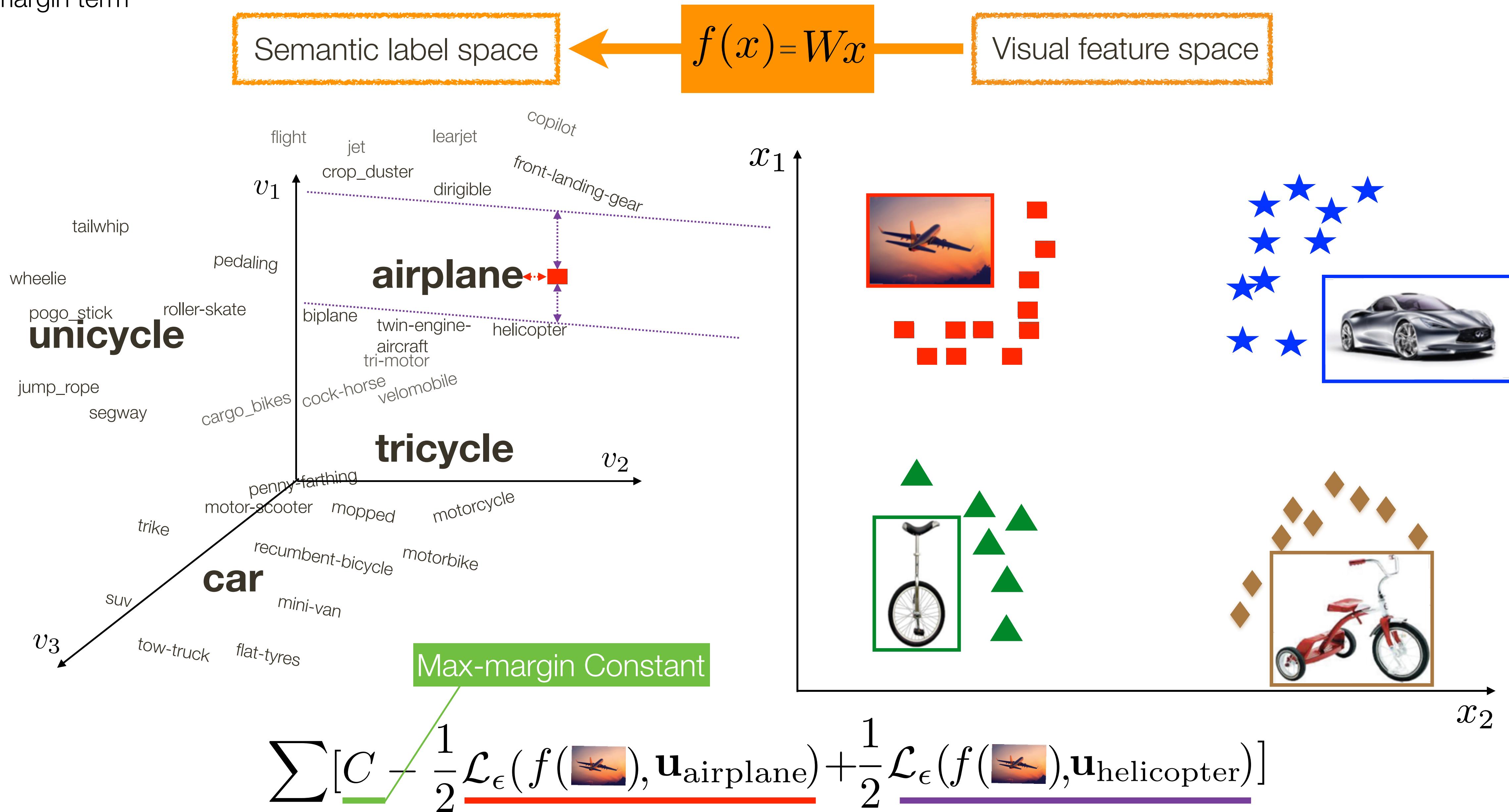
Max-margin term



$$\sum [C - \frac{1}{2} \mathcal{L}_\epsilon(f(\text{airplane}), \mathbf{u}_{\text{airplane}}) + \frac{1}{2} \mathcal{L}_\epsilon(f(\text{helicopter}), \mathbf{u}_{\text{helicopter}})]$$

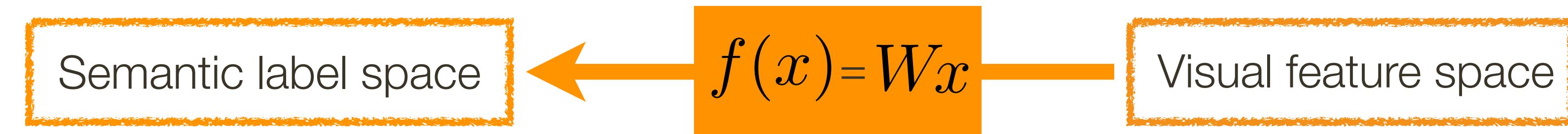
# Semi-Supervised Vocabulary-Informed Learning (SS-Voc)

Max-margin term



# Semi-Supervised Vocabulary-Informed Learning (SS-Voc)

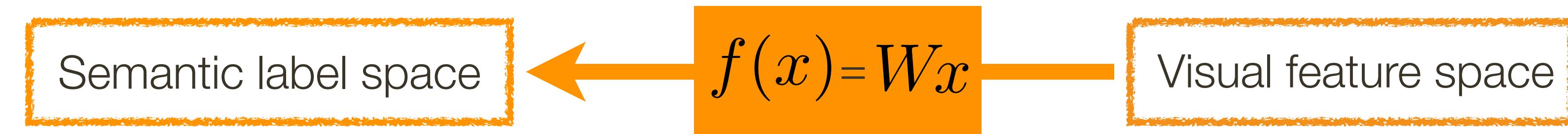
Final full objective



$$\sum \mathcal{L}_\epsilon(f(\text{airplane}, W), \mathbf{u}_{\text{airplane}} V) + \sum [C - \frac{1}{2} \mathcal{L}_\epsilon(f(\text{airplane}), \mathbf{u}_{\text{airplane}} V) + \frac{1}{2} \mathcal{L}_\epsilon(f(\text{helicopter}), \mathbf{u}_{\text{helicopter}} V)]$$

# Semi-Supervised Vocabulary-Informed Learning (SS-Voc)

Final full objective

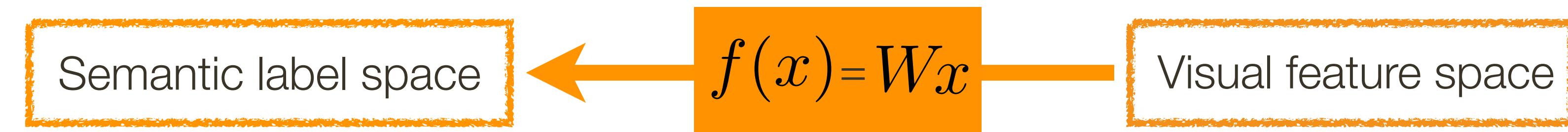


$$\sum \mathcal{L}_\epsilon(f(\text{airplane}, W), \mathbf{u}_{\text{airplane}} V) + \sum [C - \frac{1}{2} \mathcal{L}_\epsilon(f(\text{airplane}), \mathbf{u}_{\text{airplane}} V) + \frac{1}{2} \mathcal{L}_\epsilon(f(\text{helicopter}), \mathbf{u}_{\text{helicopter}} V)]$$

Regression Term

# Semi-Supervised Vocabulary-Informed Learning (SS-Voc)

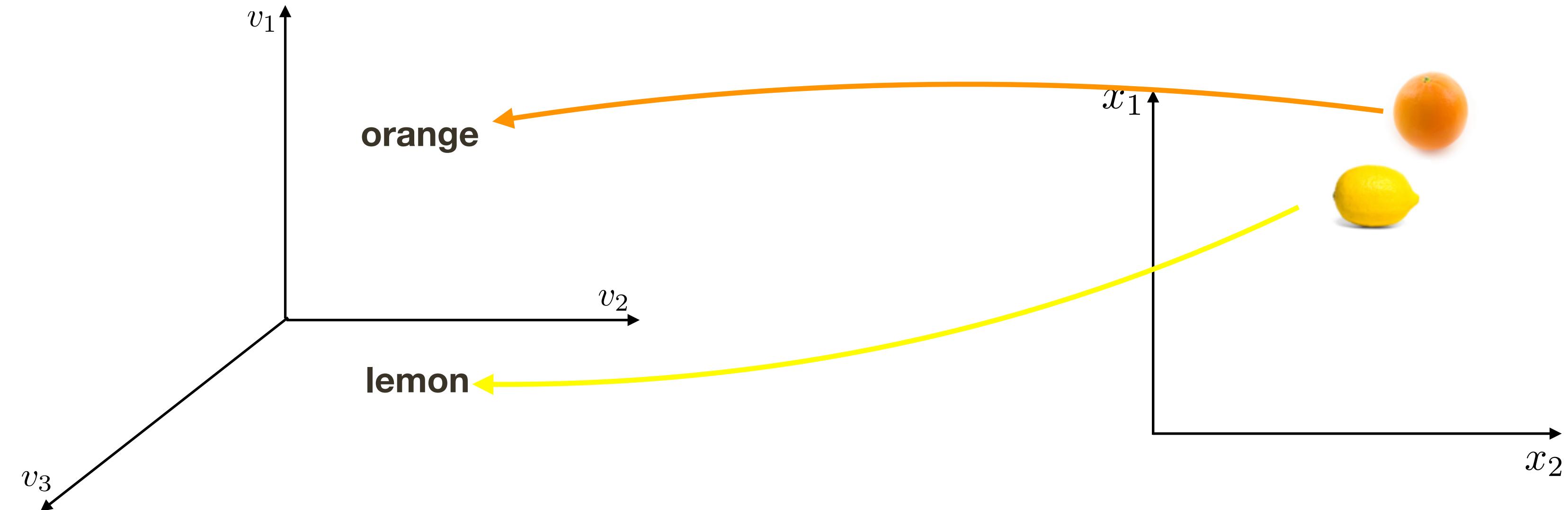
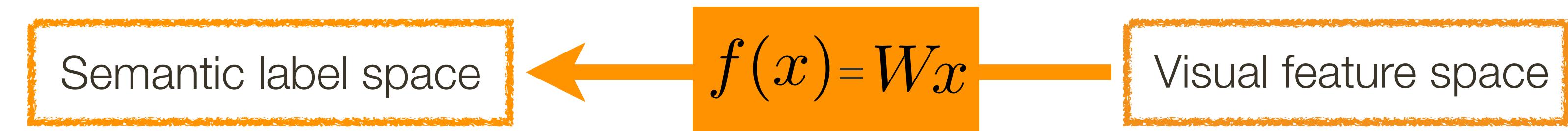
## Final full objective



$$\frac{\sum \mathcal{L}_\epsilon(f(\text{airplane}), \mathbf{u}_{\text{airplane}} V) + \sum [C - \frac{1}{2} \mathcal{L}_\epsilon(f(\text{airplane}), \mathbf{u}_{\text{airplane}} V) + \frac{1}{2} \mathcal{L}_\epsilon(f(\text{helicopter}), \mathbf{u}_{\text{helicopter}} V)]}{\text{Regression Term} \quad \quad \quad \text{Max-margin Term}}$$

# Semi-Supervised Vocabulary-Informed Learning (SS-Voc)

Final full objective



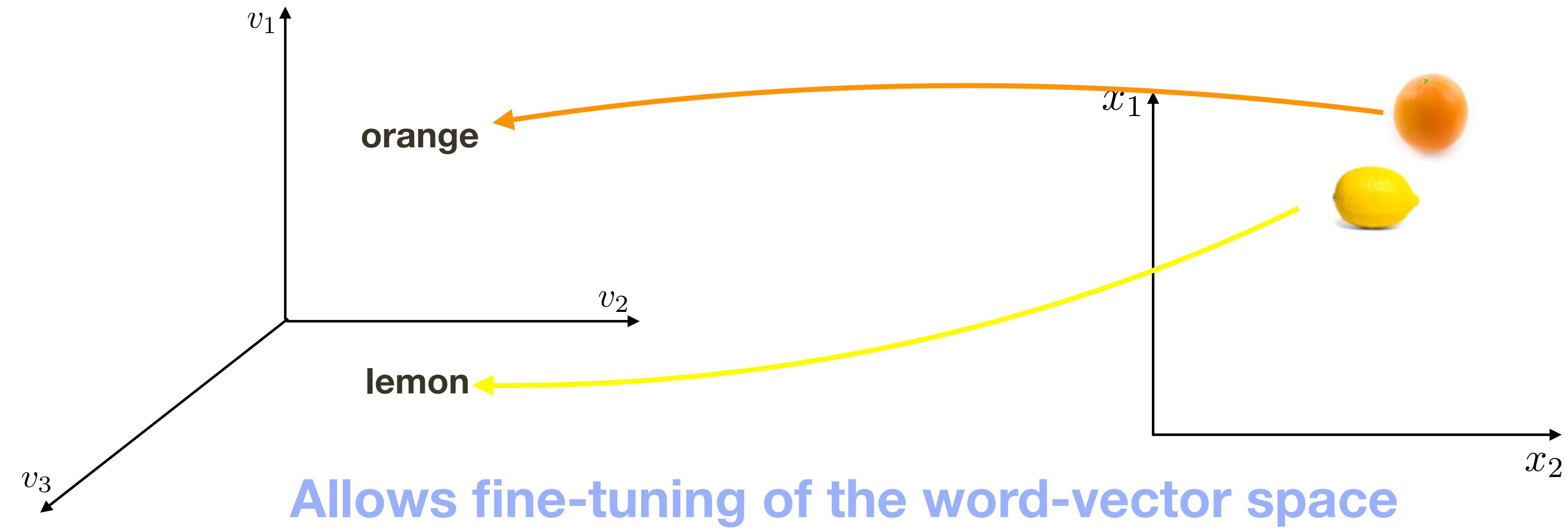
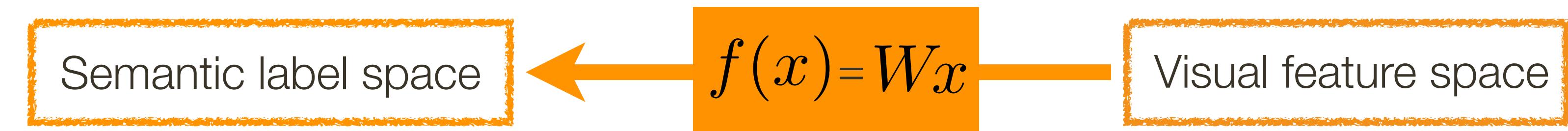
$$\sum \mathcal{L}_\epsilon(f(\text{airplane}, W), \mathbf{u}_{\text{airplane}} V) + \sum [C - \frac{1}{2} \mathcal{L}_\epsilon(f(\text{airplane}), \mathbf{u}_{\text{airplane}} V) + \frac{1}{2} \mathcal{L}_\epsilon(f(\text{helicopter}), \mathbf{u}_{\text{helicopter}} V)]$$

Regression Term

Max-margin Term

# Semi-Supervised Vocabulary-Informed Learning (SS-Voc)

Final full objective



$$\sum \mathcal{L}_\epsilon(f(\text{airplane}, W), \mathbf{u}_{\text{airplane}} | V) + \sum [C - \frac{1}{2} \mathcal{L}_\epsilon(f(\text{airplane}), \mathbf{u}_{\text{airplane}} | V) + \frac{1}{2} \mathcal{L}_\epsilon(f(\text{helicopter}), \mathbf{u}_{\text{helicopter}} | V)]$$

Word-vector space global transformation

Regression Term

Max-margin Term

# **Advantages of the Approach**

# **Advantages** of the Approach

- A new paradigm for learning informed by very large vocabulary

# Advantages of the Approach

- A new paradigm for learning informed by very large vocabulary
- A unified framework for supervised, zero-shot learning

# Advantages of the Approach

- A new paradigm for learning informed by very large vocabulary
- A unified framework for supervised, zero-shot learning
- Competitive quantitative performance

# Advantages of the Approach

- A new paradigm for learning informed by very large vocabulary
- A unified framework for supervised, zero-shot learning
- Competitive quantitative performance
- Our framework can even scale up to open set image recognition with 310,000 vocabulary entities

# Evaluation



# Datasets

# Datasets

## Animals with Attributes(AwA) [Lampert *et al.* CVPR 2009]:

40 auxiliary classes (24295 images), 10 target classes (6180 images);

We use 5 instances per auxiliary class for learning;



# Datasets

## Animals with Attributes(AwA) [Lampert *et al.* CVPR 2009]:

40 auxiliary classes (24295 images), 10 target classes (6180 images);

We use 5 instances per auxiliary class for learning;



## ImageNet 2012/2010 [Deng *et al.* CVPR 2009]:

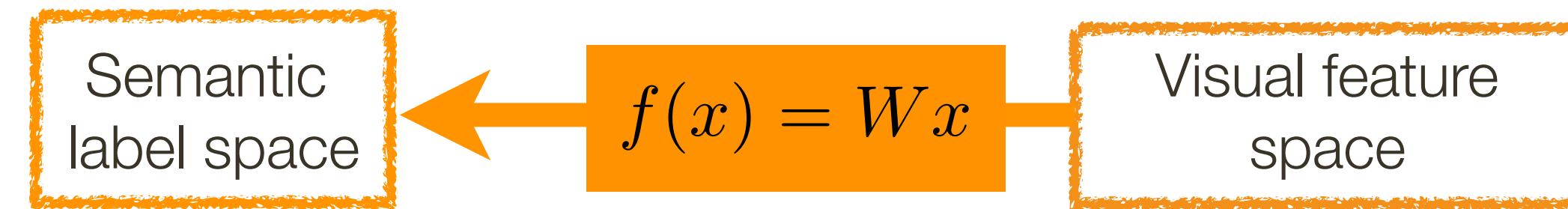
1000 auxiliary classes (from ImageNet 2012);

360 target classes (from ImageNet 2010).

We use 3 instance per auxiliary class for learning;



# Recognition Tasks

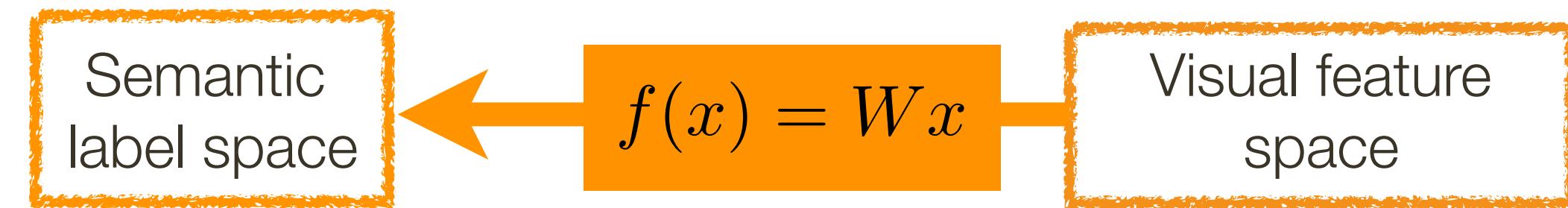


AwA/ImageNet	No. Testing Classes			No. Testing Words	
	Auxiliary	Target	Total	Vocabulary	Chance(%)
SUPERVISED	✓		40/1000	40/1000	2.5/0.1
ZERO-SHOT		✓	10/360	10/360	10/0.28
OPEN-SET	✓	✓	50/1360	310K/310K	3.2E-04

The tasks are only separated in **evaluation**;

We train **one unified SS-Voc model** for all the settings

# Recognition Tasks

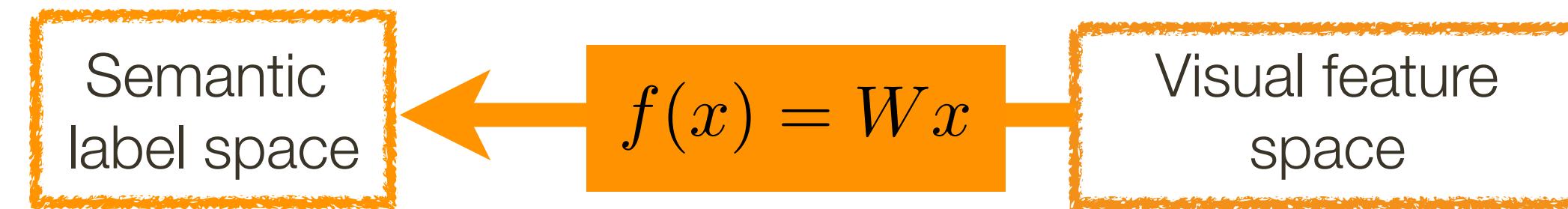


AwA/ImageNet	No. Testing Classes			No. Testing Words	
	Auxiliary	Target	Total	Vocabulary	Chance(%)
SUPERVISED	✓		40/1000	40/1000	2.5/0.1
ZERO-SHOT		✓	10/360	10/360	10/0.28
OPEN-SET	✓	✓	50/1360	310K/310K	3.2E-04

The tasks are only separated in **evaluation**;

We train **one unified SS-Voc model** for all the settings

# Recognition Tasks



AwA/ImageNet	No. Testing Classes			No. Testing Words	
	Auxiliary	Target	Total	Vocabulary	Chance(%)
SUPERVISED	✓		40/1000	40/1000	2.5/0.1
ZERO-SHOT		✓	10/360	10/360	10/0.28
OPEN-SET	✓	✓	50/1360	310K/310K	3.2E-04

The tasks are only separated in **evaluation**;

We train **one unified SS-Voc model** for all the settings

# Baselines

## SUPERVISED LEARNING

SVM: Input Features → Semantic Labels  
SVR-Map: Input Features → Semantic Word Vectors

One-shot Learning: Bart *et al.* CVPR 2005; Fei-Fei *et al.* TPAMI 2006;  
Mensink *et al.* ECCV 2012; Fu *et al.* TPAMI 2013;

# Baselines

## SUPERVISED LEARNING

SVM: Input Features → Semantic Labels  
SVR-Map: Input Features → Semantic Word Vectors

## ZERO-SHOT LEARNING

DAP/IAP, DeViSE, ConSE, AMP, PST, HEX, TMV-BLP

DAP/IAP(Lampert *et al.* TPAMI 2013); DeViSE(Frome *et al.* NIPS 2013);  
ConSE(Norouzi *et al.* ICLR 2014); AMP(Fu et al. CVPR 2015),  
PST(Rohrbach et al. NIPS 2013).

# Baselines

## SUPERVISED LEARNING

SVM: Input Features → Semantic Labels  
SVR-Map: Input Features → Semantic Word Vectors

## ZERO-SHOT LEARNING

DAP/IAP, DeViSE, ConSE, AMP, PST, HEX, TMV-BLP

## OPEN-SET IMAGE RECOGNITION

SVR-Map

Schemer *et al.* TPAMI 2013, TPAMI 2014; Sattar et al. CVPR 2015;  
Bendale et al. CVPR 2015;

# Variants of Our Model

**SS-Voc(W,V):**

$$\sum \mathcal{L}_\epsilon(f(\text{airplane}, W), \mathbf{u}_{\text{airplane}} V) + \sum [C - \frac{1}{2} \mathcal{L}_\epsilon(f(\text{airplane}), \mathbf{u}_{\text{airplane}} V) + \frac{1}{2} \mathcal{L}_\epsilon(f(\text{helicopter}), \mathbf{u}_{\text{helicopter}} V)]$$

Word-vector space global transformation

Regression Term

Max-margin Term

# Variants of Our Model

**SS-Voc(W,V):**

$$\sum \mathcal{L}_\epsilon(f(\text{airplane}, W), \mathbf{u}_{\text{airplane}}) + \sum [C - \frac{1}{2} \mathcal{L}_\epsilon(f(\text{airplane}), \mathbf{u}_{\text{airplane}}) + \frac{1}{2} \mathcal{L}_\epsilon(f(\text{helicopter}), \mathbf{u}_{\text{helicopter}})]$$

Word-vector space global transformation

Regression Term

Max-margin Term

**SS-Voc(W):**

$$\sum \mathcal{L}_\epsilon(f(\text{airplane}, W), \mathbf{u}_{\text{airplane}}) + \sum [C - \frac{1}{2} \mathcal{L}_\epsilon(f(\text{airplane}), \mathbf{u}_{\text{airplane}}) + \frac{1}{2} \mathcal{L}_\epsilon(f(\text{helicopter}), \mathbf{u}_{\text{helicopter}})]$$

Regression Term

Max-margin Term

# Variants of Our Model

**SS-Voc(W,V):**

$$\sum \mathcal{L}_\epsilon(f(\text{airplane}, W), \mathbf{u}_{\text{airplane}}) + \sum [C - \frac{1}{2} \mathcal{L}_\epsilon(f(\text{airplane}), \mathbf{u}_{\text{airplane}}) + \frac{1}{2} \mathcal{L}_\epsilon(f(\text{helicopter}), \mathbf{u}_{\text{helicopter}})]$$

Word-vector space global transformation

Regression Term

Max-margin Term

**SS-Voc(W):**

$$\sum \mathcal{L}_\epsilon(f(\text{airplane}, W), \mathbf{u}_{\text{airplane}}) + \sum [C - \frac{1}{2} \mathcal{L}_\epsilon(f(\text{airplane}), \mathbf{u}_{\text{airplane}}) + \frac{1}{2} \mathcal{L}_\epsilon(f(\text{helicopter}), \mathbf{u}_{\text{helicopter}})]$$

Regression Term

Max-margin Term

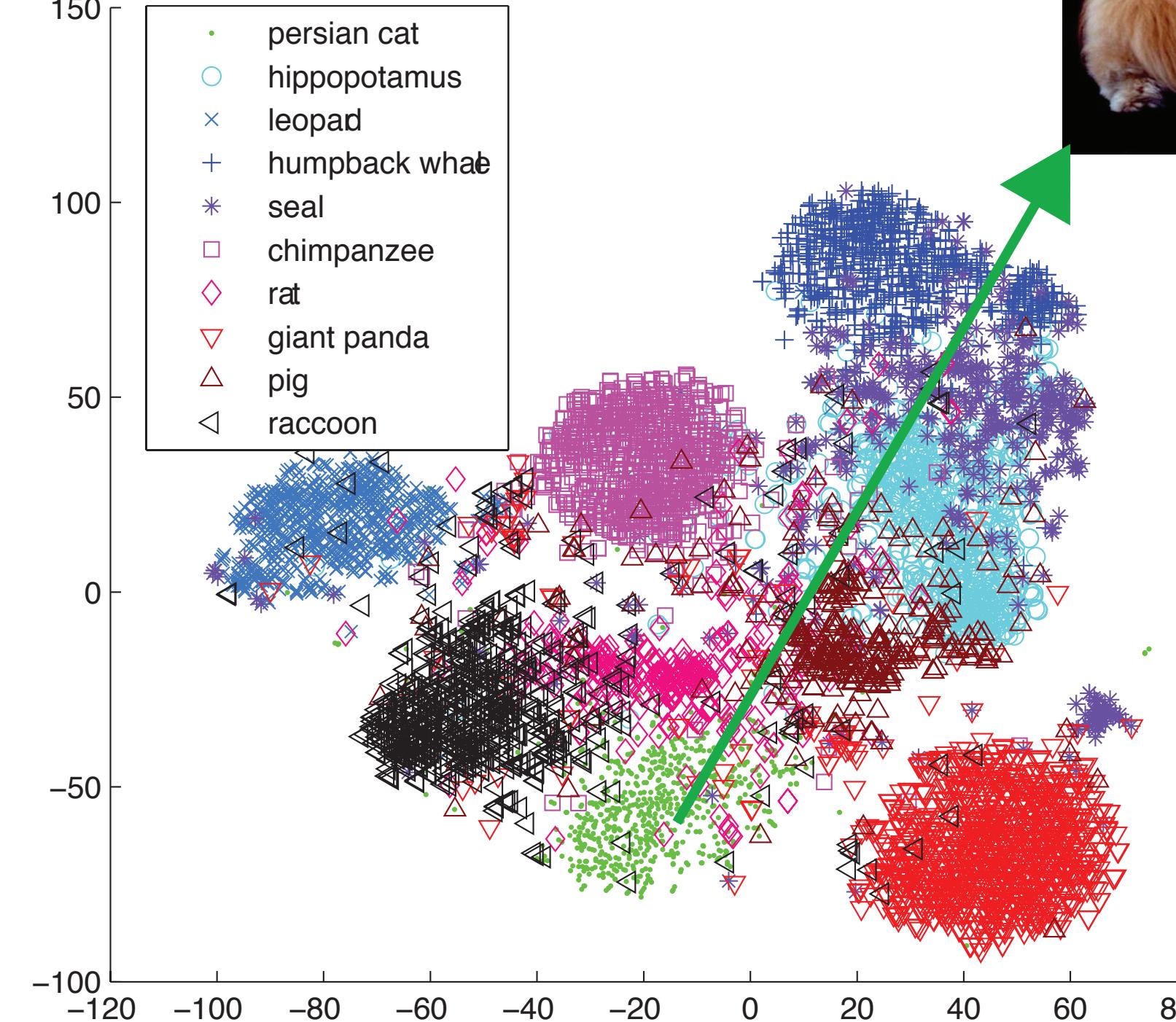
**SVR-Map:**  $\sum \mathcal{L}_\epsilon(f(\text{airplane}, W), \mathbf{u}_{\text{airplane}})$

Regression Term

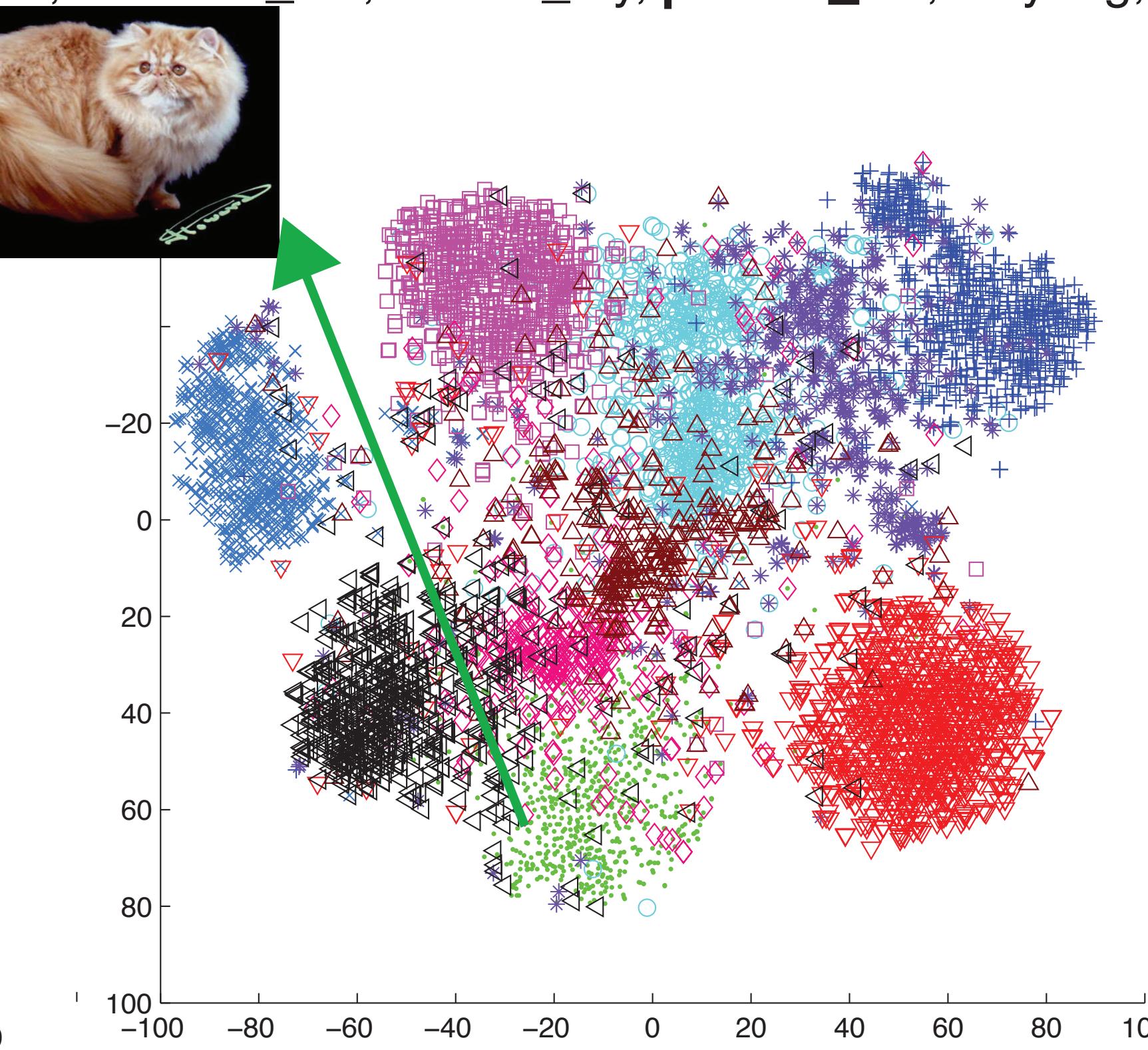
# t-SNE Visualization of AwA 10 Testing Classes

SS-Voc: **persian\_cat**, siamese\_cat, hamster, weasel, rabbit, monkey, zebra, owl, anthropomorphized, cat

SVR-Map: hamster, squirrel, rabbit, raccoon, kitten, siamese\_cat, stuffed\_toy, **persian\_cat**, ladybug, puppy



**Semi-Supervised Vocabulary-informed  
Learning (SS-Voc)**



**Support Vector Regression (SVR)**

# Supervised Results

## AwA dataset

Method	Accuracy
SS-Voc (W,V)	59.1
SS-Voc (W)	58.6
SVM	52.1
SVR-Map	57.1

## ImageNet dataset

Method	Accuracy
SS-Voc (W,V)	37.1
SS-Voc (W)	36.3
SVM	33.8
SVR-Map	25.6

# Supervised Results

## AwA dataset

Method	Accuracy
SS-Voc (W,V)	59.1
SS-Voc (W)	58.6
SVM	52.1
SVR-Map	57.1

## ImageNet dataset

Method	Accuracy
SS-Voc (W,V)	37.1
SS-Voc (W)	36.3
SVM	33.8
SVR-Map	25.6

# Supervised Results

## AwA dataset

Method	Accuracy
SS-Voc (W,V)	59.1
SS-Voc (W)	58.6
SVM	52.1
SVR-Map	57.1

## ImageNet dataset

Method	Accuracy
SS-Voc (W,V)	37.1
SS-Voc (W)	36.3
SVM	33.8
SVR-Map	25.6

# Supervised Results

## AwA dataset

Method	Accuracy
SS-Voc (W,V)	59.1
SS-Voc (W)	58.6
SVM	52.1
SVR-Map	57.1

## ImageNet dataset

Method	Accuracy
SS-Voc (W,V)	37.1
SS-Voc (W)	36.3
SVM	33.8
SVR-Map	25.6

# Supervised Results

## AwA dataset

Method	Accuracy
SS-Voc (W,V)	59.1
SS-Voc (W)	58.6
SVM	52.1
SVR-Map	57.1

## ImageNet dataset

Method	Accuracy
SS-Voc (W,V)	37.1
SS-Voc (W)	36.3
SVM	33.8
SVR-Map	25.6

# Supervised Results

**AwA dataset**

Method	Accuracy
SS-Voc (W,V)	59.1
SS-Voc (W)	58.6
SVM	52.1
SVR-Map	57.1

**ImageNet dataset**

Method	Accuracy
SS-Voc (W,V)	37.1
SS-Voc (W)	36.3
SVM	33.8
SVR-Map	25.6

# Supervised Results

## AwA dataset

Method	Accuracy
SS-Voc (W,V)	59.1
SS-Voc (W)	58.6
SVM	52.1
SVR-Map	57.1

## ImageNet dataset

Method	Accuracy
SS-Voc (W,V)	37.1
SS-Voc (W)	36.3
SVM	33.8
SVR-Map	25.6

# Supervised Results

**AwA dataset**

Method	Accuracy
SS-Voc (W,V)	59.1
SS-Voc (W)	58.6
SVM	52.1
SVR-Map	57.1

**ImageNet dataset**

Method	Accuracy
SS-Voc (W,V)	37.1
SS-Voc (W)	36.3
SVM	33.8
SVR-Map	25.6

# Supervised Results

## AwA dataset

Method	Accuracy
SS-Voc (W,V)	59.1
SS-Voc (W)	58.6
SVM	52.1
SVR-Map	57.1

## ImageNet dataset

Method	Accuracy
SS-Voc (W,V)	37.1
SS-Voc (W)	36.3
SVM	33.8
SVR-Map	25.6

# Zero-shot Results—AwA dataset

Method	Features	Accuracy
SS-Voc: full instances  800 instances (20 inst*40 class);  200 instances (5 inst*40 class);	CNN <sub>OverFeat</sub>	78.3
	CNN <sub>OverFeat</sub>	74.4
	CNN <sub>OverFeat</sub>	68.9
Akata <i>et al.</i> CVPR 2015	CNN <sub>GoogLeNet</sub>	73.9
TMV-BLP (Fu <i>et al.</i> ECCV 2014)	CNN <sub>OverFeat</sub>	69.9
AMP (SR+SE) (Fu <i>et al.</i> CVPR 2015)	CNN <sub>OverFeat</sub>	66.0
DAP (Lampert <i>et al.</i> TPAMI 2013)	CNN <sub>VGG19</sub>	57.5
PST (Rohrbach <i>et al.</i> NIPS 2013)	CNN <sub>OverFeat</sub>	53.2
DS (Rohrbach <i>et al.</i> CVPR 2010)	CNN <sub>OverFeat</sub>	52.7
IAP (Lampert <i>et al.</i> TPAMI 2013)	CNN <sub>OverFeat</sub>	44.5
HEX (Deng <i>et al.</i> ECCV 2014)	CNN <sub>DECAF</sub>	44.2

# Zero-shot Results—AwA dataset

Method	Features	Accuracy
SS-Voc: full instances	CNN <sub>OverFeat</sub>	78.3
800 instances (20 inst*40 class); 200 instances (5 inst*40 class);	CNN <sub>OverFeat</sub>	74.4
	CNN <sub>OverFeat</sub>	68.9
Akata <i>et al.</i> CVPR 2015	CNN <sub>GoogLeNet</sub>	73.9
TMV-BLP (Fu <i>et al.</i> ECCV 2014)	CNN <sub>OverFeat</sub>	69.9
AMP (SR+SE) (Fu <i>et al.</i> CVPR 2015)	CNN <sub>OverFeat</sub>	66.0
DAP (Lampert <i>et al.</i> TPAMI 2013)	CNN <sub>VGG19</sub>	57.5
PST (Rohrbach <i>et al.</i> NIPS 2013)	CNN <sub>OverFeat</sub>	53.2
DS (Rohrbach <i>et al.</i> CVPR 2010)	CNN <sub>OverFeat</sub>	52.7
IAP (Lampert <i>et al.</i> TPAMI 2013)	CNN <sub>OverFeat</sub>	44.5
HEX (Deng <i>et al.</i> ECCV 2014)	CNN <sub>DECAF</sub>	44.2

# Zero-shot Results—AwA dataset

Method	Features	Accuracy
SS-Voc: full instances  800 instances (20 inst*40 class);  200 instances (5 inst*40 class);	CNN <sub>OverFeat</sub>	78.3
	CNN <sub>OverFeat</sub>	74.4
	CNN <sub>OverFeat</sub>	68.9
Akata <i>et al.</i> CVPR 2015	CNN <sub>GoogLeNet</sub>	73.9
TMV-BLP (Fu <i>et al.</i> ECCV 2014)	CNN <sub>OverFeat</sub>	69.9
AMP (SR+SE) (Fu <i>et al.</i> CVPR 2015)	CNN <sub>OverFeat</sub>	66.0
DAP (Lampert <i>et al.</i> TPAMI 2013)	CNN <sub>VGG19</sub>	57.5
PST (Rohrbach <i>et al.</i> NIPS 2013)	CNN <sub>OverFeat</sub>	53.2
DS (Rohrbach <i>et al.</i> CVPR 2010)	CNN <sub>OverFeat</sub>	52.7
IAP (Lampert <i>et al.</i> TPAMI 2013)	CNN <sub>OverFeat</sub>	44.5
HEX (Deng <i>et al.</i> ECCV 2014)	CNN <sub>DECAF</sub>	44.2

# Zero-shot Results—AwA dataset

Method	Features	Accuracy
SS-Voc: full instances	CNN <sub>OverFeat</sub>	78.3
800 instances (20 inst*40 class); 200 instances (5 inst*40 class); Akata <i>et al.</i> CVPR 2015	CNN <sub>OverFeat</sub>	74.4
TMV-BLP (Fu <i>et al.</i> ECCV 2014)	CNN <sub>OverFeat</sub>	68.9
AMP (SR+SE) (Fu <i>et al.</i> CVPR 2015)	CNN <sub>OverFeat</sub>	73.9
DAP (Lampert <i>et al.</i> TPAMI 2013)	CNN <sub>VGG19</sub>	69.9
PST (Rohrbach <i>et al.</i> NIPS 2013)	CNN <sub>OverFeat</sub>	66.0
DS (Rohrbach <i>et al.</i> CVPR 2010)	CNN <sub>OverFeat</sub>	57.5
IAP (Lampert <i>et al.</i> TPAMI 2013)	CNN <sub>OverFeat</sub>	53.2
IAP (Lampert <i>et al.</i> TPAMI 2013)	CNN <sub>OverFeat</sub>	52.7
HEX (Deng <i>et al.</i> ECCV 2014)	CNN <sub>DECAF</sub>	44.5
HEX (Deng <i>et al.</i> ECCV 2014)	CNN <sub>DECAF</sub>	44.2

3.3%

SS-Voc: full instances

# Zero-shot Results—AwA dataset

Method	Features	Accuracy
SS-Voc: full instances  800 instances (20 inst*40 class);  200 instances (5 inst*40 class);	CNN <sub>OverFeat</sub>	78.3
	CNN <sub>OverFeat</sub>	74.4
	CNN <sub>OverFeat</sub>	68.9
Akata <i>et al.</i> CVPR 2015	CNN <sub>GoogLeNet</sub>	73.9
TMV-BLP (Fu <i>et al.</i> ECCV 2014)	CNN <sub>OverFeat</sub>	69.9
AMP (SR+SE) (Fu <i>et al.</i> CVPR 2015)	CNN <sub>OverFeat</sub>	66.0
DAP (Lampert <i>et al.</i> TPAMI 2013)	CNN <sub>VGG19</sub>	57.5
PST (Rohrbach <i>et al.</i> NIPS 2013)	CNN <sub>OverFeat</sub>	53.2
DS (Rohrbach <i>et al.</i> CVPR 2010)	CNN <sub>OverFeat</sub>	52.7
IAP (Lampert <i>et al.</i> TPAMI 2013)	CNN <sub>OverFeat</sub>	44.5
HEX (Deng <i>et al.</i> ECCV 2014)	CNN <sub>DECAF</sub>	44.2

# Zero-shot Results—AwA dataset

0.82%

Method	Features	Accuracy
SS-Voc: full instances	CNN <sub>OverFeat</sub>	78.3
800 instances (20 inst*40 class);	CNN <sub>OverFeat</sub>	74.4
200 instances (5 inst*40 class);	CNN <sub>OverFeat</sub>	68.9
Akata <i>et al.</i> CVPR 2015	CNN <sub>GoogLeNet</sub>	73.9
TMV-BLP (Fu <i>et al.</i> ECCV 2014)	CNN <sub>OverFeat</sub>	69.9
AMP (SR+SE) (Fu <i>et al.</i> CVPR 2015)	CNN <sub>OverFeat</sub>	66.0
DAP (Lampert <i>et al.</i> TPAMI 2013)	CNN <sub>VGG19</sub>	57.5
PST (Rohrbach <i>et al.</i> NIPS 2013)	CNN <sub>OverFeat</sub>	53.2
DS (Rohrbach <i>et al.</i> CVPR 2010)	CNN <sub>OverFeat</sub>	52.7
IAP (Lampert <i>et al.</i> TPAMI 2013)	CNN <sub>OverFeat</sub>	44.5
HEX (Deng <i>et al.</i> ECCV 2014)	CNN <sub>DECAF</sub>	44.2

# Zero-shot Results—AwA dataset

Method	Features	Accuracy
SS-Voc: full instances  800 instances (20 inst*40 class);  200 instances (5 inst*40 class);	CNN <sub>OverFeat</sub>	78.3
	CNN <sub>OverFeat</sub>	74.4
	CNN <sub>OverFeat</sub>	68.9
Akata <i>et al.</i> CVPR 2015	CNN <sub>GoogLeNet</sub>	73.9
TMV-BLP (Fu <i>et al.</i> ECCV 2014)	CNN <sub>OverFeat</sub>	69.9
AMP (SR+SE) (Fu <i>et al.</i> CVPR 2015)	CNN <sub>OverFeat</sub>	66.0
DAP (Lampert <i>et al.</i> TPAMI 2013)	CNN <sub>VGG19</sub>	57.5
PST (Rohrbach <i>et al.</i> NIPS 2013)	CNN <sub>OverFeat</sub>	53.2
DS (Rohrbach <i>et al.</i> CVPR 2010)	CNN <sub>OverFeat</sub>	52.7
IAP (Lampert <i>et al.</i> TPAMI 2013)	CNN <sub>OverFeat</sub>	44.5
HEX (Deng <i>et al.</i> ECCV 2014)	CNN <sub>DECAF</sub>	44.2

# Zero-shot Results—ImageNet

Method	Features	T-1 Accuracy (full instances)	T-5 Accuracy (full instances)	T-1 Accuracy (3000 instances)	T-5 Accuracy (3000 instances)
SS-Voc	CNN <sub>VGG-19</sub>	9.5	16.8	8.9	14.9
ConSE	CNN <sub>VGG-19</sub>	7.8	15.5	5.5	13.1
DeViSE	CNN <sub>VGG-19</sub>	5.2	12.8	3.7	11.8
AMP	CNN <sub>VGG-19</sub>	6.1	13.1	3.5	10.5

# Zero-shot Results—ImageNet

Method	Features	T-1 Accuracy (full instances)	T-5 Accuracy (full instances)	T-1 Accuracy (3000 instances)	T-5 Accuracy (3000 instances)
SS-Voc	CNN <sub>VGG-19</sub>	9.5	16.8	8.9	14.9
ConSE	CNN <sub>VGG-19</sub>	7.8	15.5	5.5	13.1
DeViSE	CNN <sub>VGG-19</sub>	5.2	12.8	3.7	11.8
AMP	CNN <sub>VGG-19</sub>	6.1	13.1	3.5	10.5

# Zero-shot Results—ImageNet

Method	Features	T-1 Accuracy (full instances)	T-5 Accuracy (full instances)	T-1 Accuracy (3000 instances)	T-5 Accuracy (3000 instances)
SS-Voc	CNN <sub>VGG-19</sub>	9.5	16.8	8.9	14.9
ConSE	CNN <sub>VGG-19</sub>	7.8	15.5	5.5	13.1
DeViSE	CNN <sub>VGG-19</sub>	5.2	12.8	3.7	11.8
AMP	CNN <sub>VGG-19</sub>	6.1	13.1	3.5	10.5

# Open-Set Image Recognition – AwA dataset

AwA	Testing Classes			Vocabulary
	Auxiliary	Target	Total	
OPEN-SET <sub>1K-NN</sub>			40/10	1000*
OPEN-SET <sub>1K-RND</sub>	(LEFT)	(RIGHT)	40/10	1000†
OPEN-SET <sub>310K</sub>			40/10	310K

# Open-Set Image Recognition – AwA dataset

AwA	Testing Classes			Vocabulary
	Auxiliary	Target	Total	
OPEN-SET <sub>1K-NN</sub>			40/10	1000*
OPEN-SET <sub>1K-RND</sub>	(LEFT)	(RIGHT)	40/10	1000†
OPEN-SET <sub>310K</sub>			40/10	310K

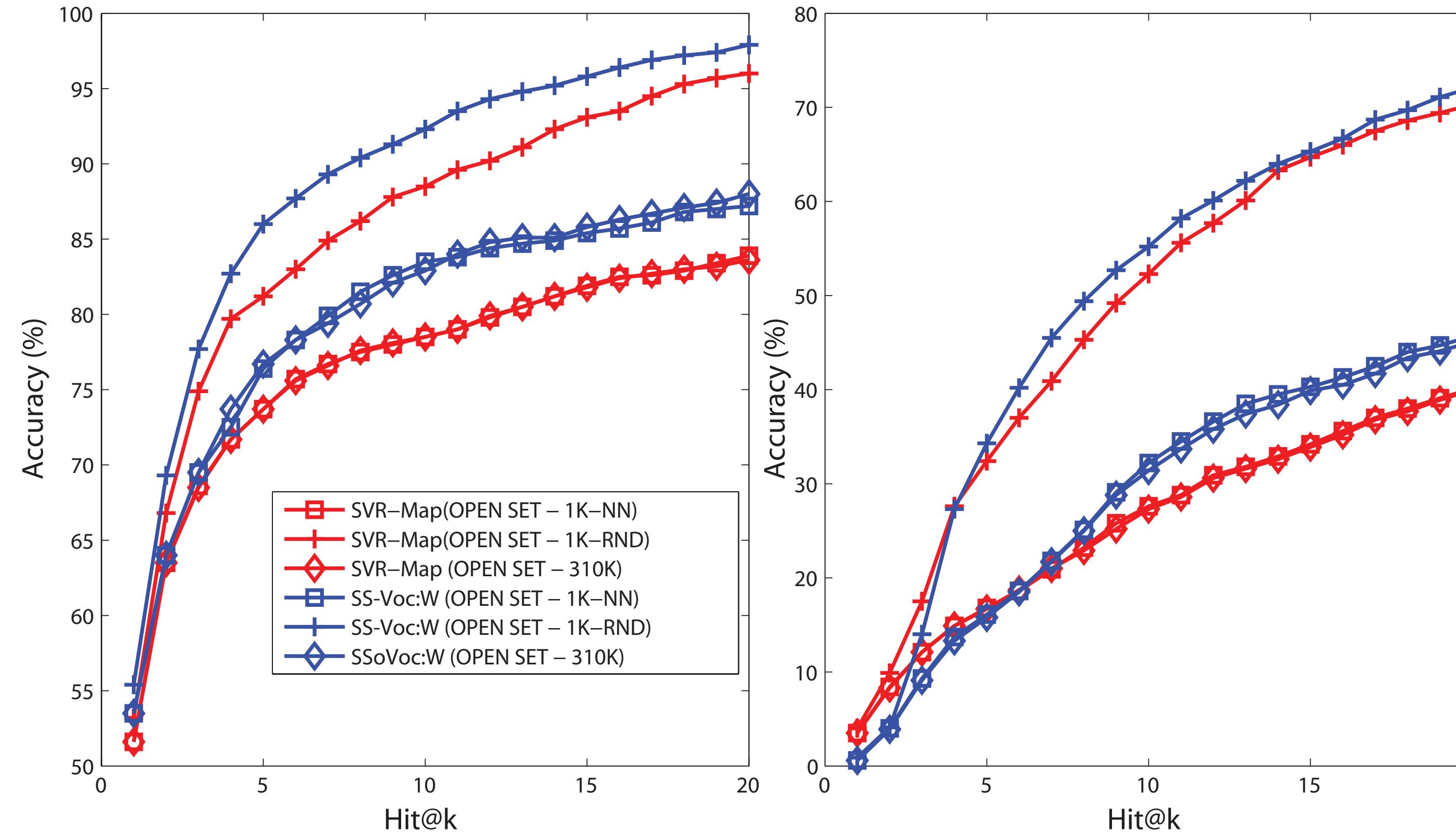
OPEN-SET <sub>1K-NN</sub>	1000 candidate labels (of 310K labels) sampled from nearest neighbor set of ground-truth class prototypes
OPEN-SET <sub>1K-RND</sub>	1000 candidate labels randomly sampled from 310K vocabulary set.
OPEN-SET <sub>310K</sub>	the large vocabulary of approximately 310K entities

# Open-Set Image Recognition – AwA dataset

AwA	Testing Classes			Vocabulary
	Auxiliary	Target	Total	
OPEN-SET <sub>1K-NN</sub>			40/10	1000*
OPEN-SET <sub>1K-RND</sub>	(LEFT)	(RIGHT)	40/10	1000†
OPEN-SET <sub>310K</sub>			40/10	310K

# Open-Set Image Recognition – AwA dataset

AwA	Testing Classes			Vocabulary
	Auxiliary	Target	Total	
OPEN-SET <sub>1K-NN</sub>			40/10	1000*
OPEN-SET <sub>1K-RND</sub>	(LEFT)(RIGHT)		40/10	1000†
OPEN-SET <sub>310K</sub>			40/10	310K



# Take-home

1. A new learning paradigm – vocabulary-informed learning
2. A unified semantic embedding framework for supervised, zero-shot and open-set image recognition

# Thanks! Q&A

