

Marcin Marszałek (marszale@inrialpes.fr)
INRIA Grenoble, LEAR / LJK

Ivan Laptev (ivan.laptev@inria.fr)
INRIA Rennes, IRISA

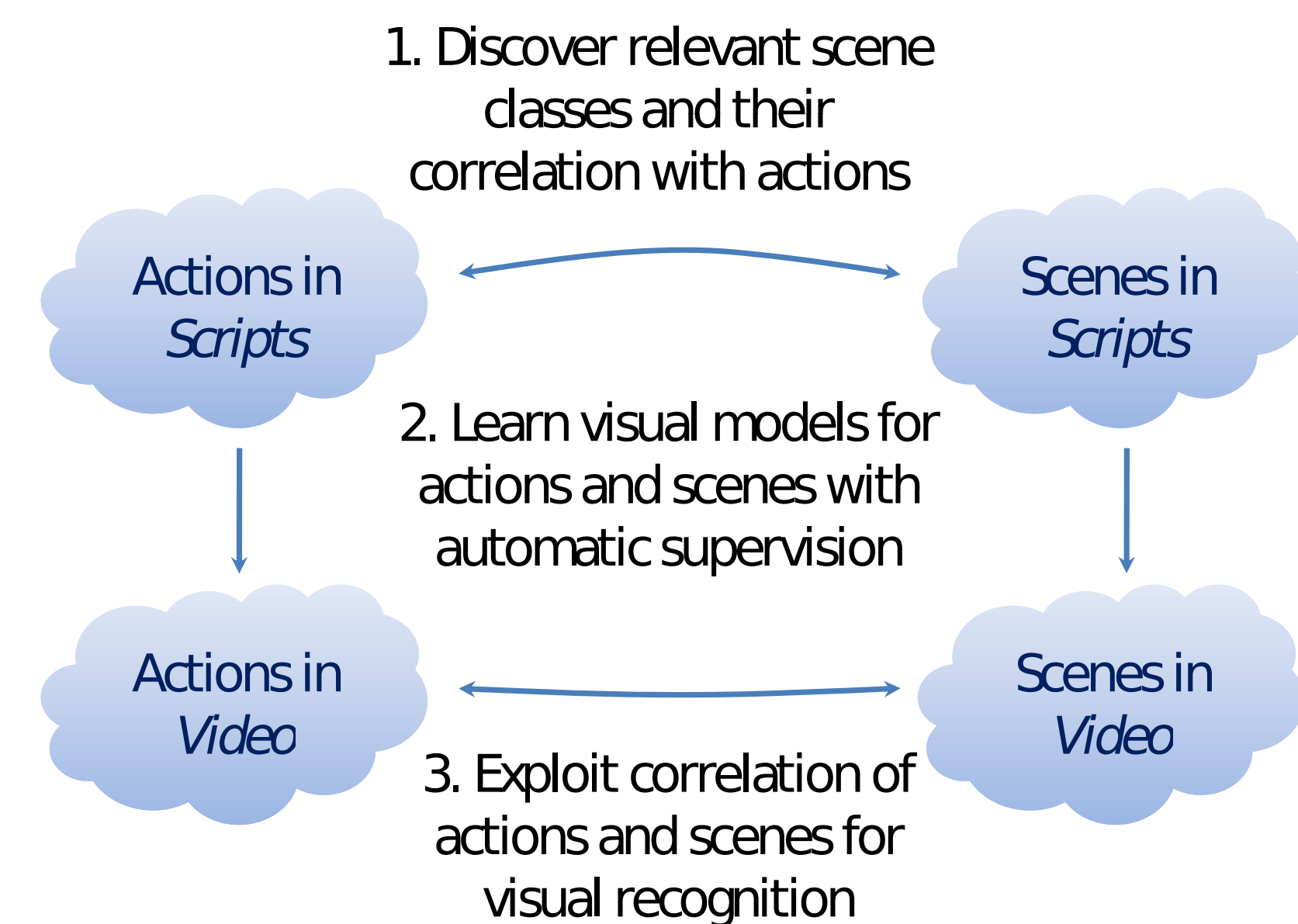
Cordelia Schmid (schmid@inrialpes.fr)
INRIA Grenoble, LEAR / LJK

Motivation and Approach

Human actions are frequently correlated with particular scene classes due to *functional* and *physical* properties of the scenes:



Moreover, some actions are *defined* by the scene context:



Movie Script Mining

We use *movie scripts* aligned with videos to:

- Discover co-occurrence relations between actions and scenes
- Automatically collect video samples for training

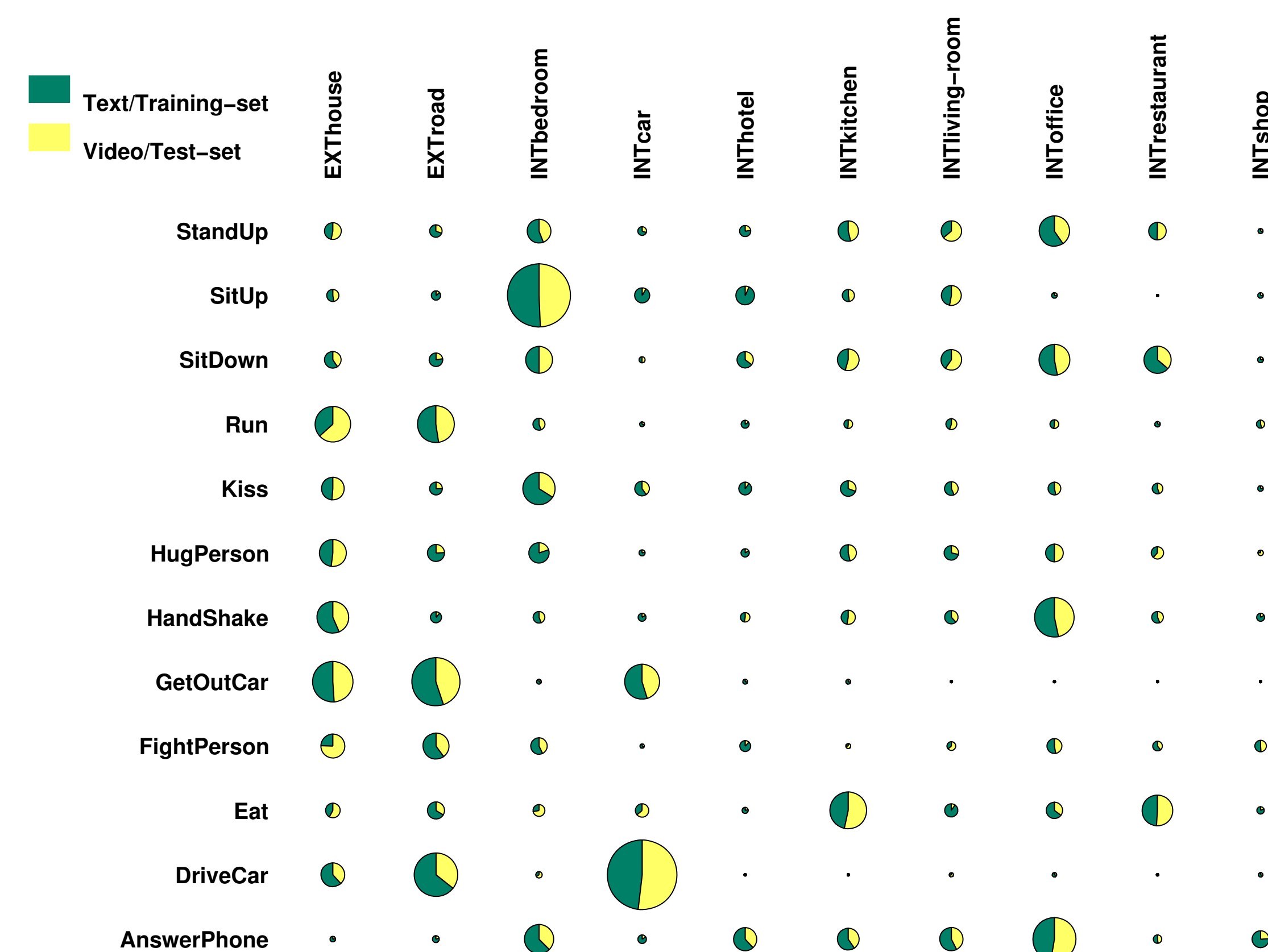
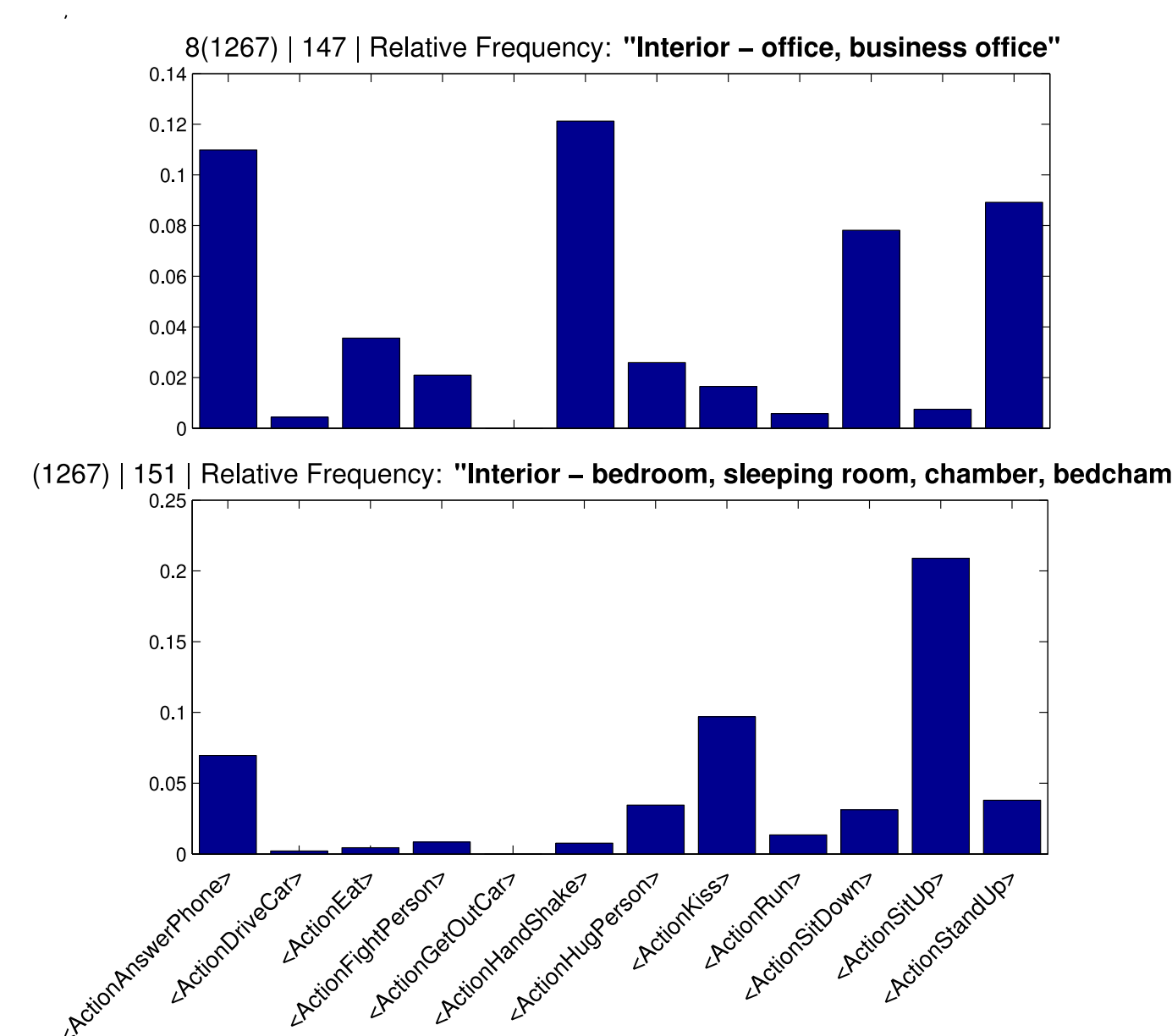
Dataset:

- Video samples are obtained from 33 training and 36 test movies
- 12 action classes are distributed among 810 automatically generated training samples and 884 manually verified test samples (approx. 7 hours of video in total)
- 10 scene classes are distributed among 570 automatically generated training samples and 582 manually verified test samples (approx. 11 hours of video in total)
- Actions-scenes co-occurrence is estimated from a large independent set of movie scripts
- The dataset is available from <http://www.irisa.fr/vista/actions/hollywood2>

Subtitles	Scene caption	Script
00:24:22 → 00:24:25 — Yes, Monsieur Laszlo. Right this way.	int. Rick's cafe, main room, night	Speech Monsieur Laszlo. Right this way. Scene description As the headwaiter takes them to a table they pass by the piano, and the woman looks at Sam. Sam, with a conscious effort, keeps his eyes on the keyboard as they go past. The headwaiter seats Ilsa...
00:24:51 → 00:24:53 Two Cointreaux, please.		Speech Two cointreaux, please.

Procedure:

- Label action samples in scripts using action text classifier
- Find frequent words and word pairs in scene captions
- Perform semantic stemming using WordNet
- Select words with high co-occurrence w.r.t. given actions
- Re-order words by the entropy $S(x)$, $x = p(action|word)$



$p(Scene|Action)$ estimated from scripts (green) and ground truth visual annotation (yellow). Note that the discovered correlations are not only intuitive, but also consistent between text and vision.

Visual Learning

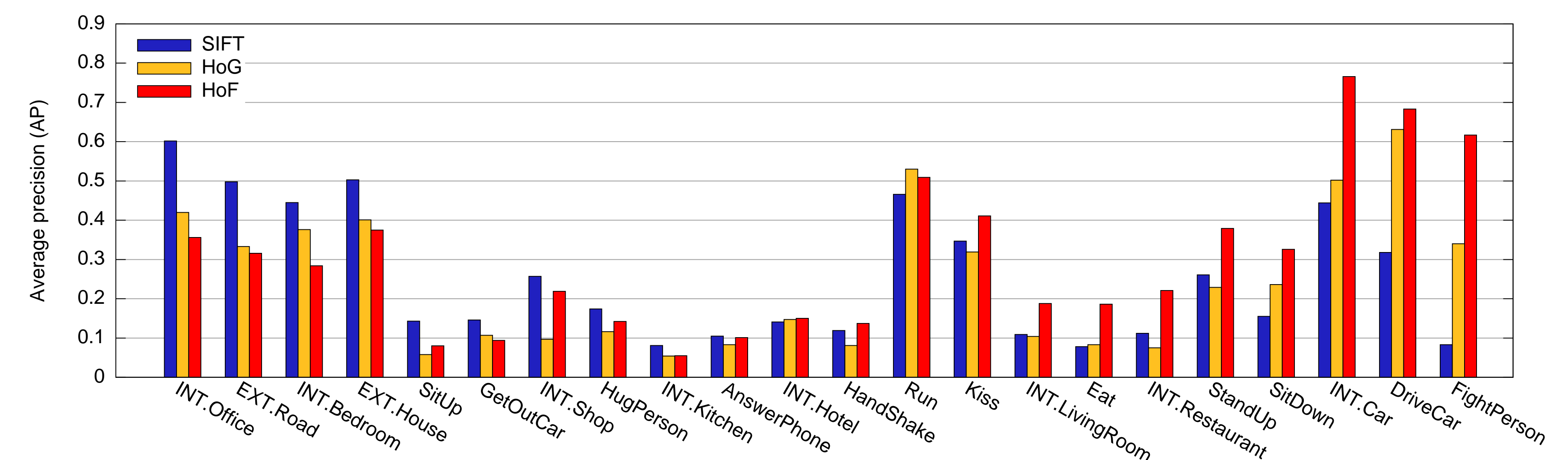
Interest points for a movie frame. 3D Harris (left) focuses on motion, whereas 2D Harris (right) regions are distributed over the scene.



We use

- Combination of local static and dynamic features:
 - 2D Harris detector + SIFT descriptor (static appearance)
 - 3D Harris detector + space-time HOG descr. (dynamic appearance)
 - 3D Harris detector + space-time HOF descriptor (motion)
- Video representation by histograms of quantized local features
- SVMs with χ^2 kernel for classification

Comparison of single feature types using bag-of-features classification approach. The static SIFT features perform well for most scene types while space-time HoG and HoF features dominate for action recognition



Classification with Context

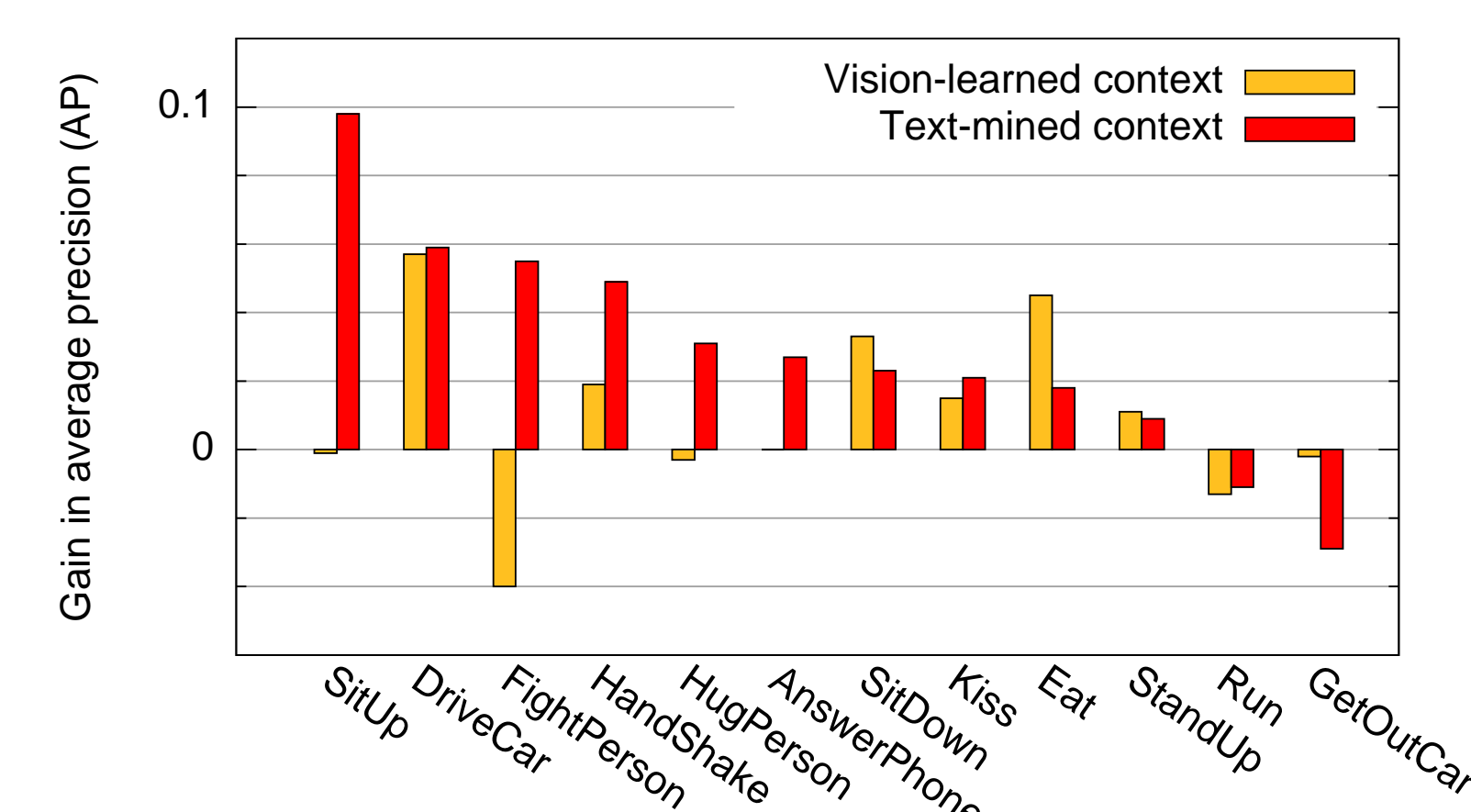
We integrate context by updating the classification score $g_a(x)$ for an action $a \in \mathcal{A}$ with a linear combination of context scores $g_s(x)$ for scene classes $s \in \mathcal{S}$:

$$g'_a(x) = g_a(x) + \tau \sum_{s \in \mathcal{S}} w_{as} g_s(x)$$

where τ is a global context weight and w_{as} are weights linking concepts a and s . We explore two ways to obtain w_{as} :

- from text – we set $w_{as} = p(s \in \mathcal{S} | a \in \mathcal{A})$
- from visual data – we train a second-layer linear SVM

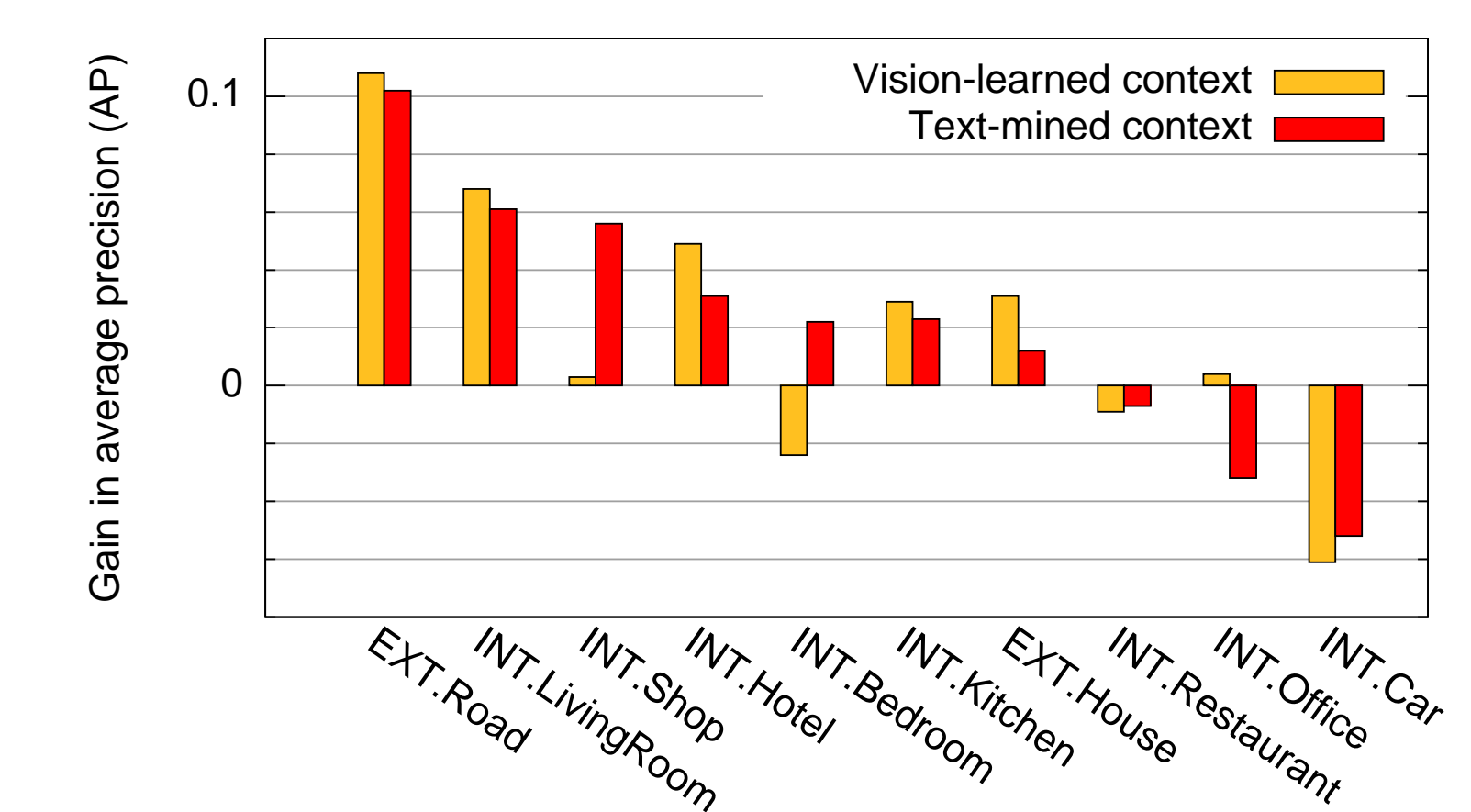
Exploiting scene context in action recognition. Note the consistent improvement for most action classes



Mean Average Precision (MAP) for action and scene classification with and without context. We also compare to chance level and try context only.

Actions	MAP	Scenes	MAP
text context	0.355	text context	0.373
vision context	0.336	vision context	0.371
no context	0.325	no context	0.351
context only	0.238	context only	0.277
chance	0.125	chance	0.162

Exploiting action context in scene recognition. Note the significant improvement for the leftmost categories



Test samples of actions where scene context significantly helps action recognition

