

Chapter 2

An Approach to Interactive Co-segmentation

Abstract In this chapter, we describe in detail our approach to interactive co-segmentation. We formulate the task as an energy minimization problem across all related images in a group. The energies across images are tied together via a shared appearance model, thus allowing for efficient inference. After describing our formulation, we present an active learning approach that makes efficient use of users' time. A wide variety of cues are combined to intelligently guide the users' next scribbles. We then introduce our co-segmentation dataset, The CMU-Cornell iCoseg dataset, the largest of its kind to date. We evaluate our system on this dataset using machine simulations as well as real user-studies. We find that our approach can achieve comparable co-segmentation performance with less user effort.

2.1 Energy Minimization

Our approach [3] to multiple-image interactive co-segmentation is a natural extension of [7]. Given user scribbles indicating foreground / background, we cast our labelling problem as minimization of Gibbs energies defined over graphs constructed over each image in a group. Specifically, consider a group of m image-scribble pairs $D = \{(\mathcal{X}^{(1)}, \mathcal{S}^{(1)}), \dots, (\mathcal{X}^{(m)}, \mathcal{S}^{(m)})\}$, where the k^{th} image is represented as a collection of n_k sites to be labelled, *i.e.* $\mathcal{X}^{(k)} = \{X_1^{(k)}, X_2^{(k)}, \dots, X_{n_k}^{(k)}\}$, and scribbles for an image $\mathcal{S}^{(k)}$ are represented as the partial (potentially empty)¹ set of labels for these sites. For computational efficiency, we use superpixels as these labelling sites (instead of pixels).² For each image (k), we build a graph, $\mathcal{G}^{(k)} = (\mathcal{V}^{(k)}, \mathcal{E}^{(k)})$, over superpixels, with edges connecting adjacent superpixels.

¹ Specifically, we require at least one labelled foreground and background site to train our models, but only one per *group*, not per image.

² We use mean-shift [10] to extract these superpixels, and typically break down 350×500 images into 400 superpixels per image.

At the start of our algorithm, we require at least one foreground and background scribble each. They can be in the same image, or in multiple images. Subsequent iterations can have a scribble just from foreground or background. Using these labelled sites, we learn a group appearance model $\mathcal{A} = \{A_1, A_2\}$, where A_1 is the first-order (unary) appearance model, and A_2 the second-order (pairwise) appearance model. This appearance model (\mathcal{A}) is described in detail in the following sections. We note that all images in the group share a common model, *i.e.* only one model is learnt. Using this appearance model, we define a collection of energies over each of the m images as follows:

$$E^{(k)}(\mathcal{X}^{(k)} : \mathcal{A}) = \sum_{i \in \mathcal{V}^{(k)}} E_i(X_i^{(k)} : A_1) + \lambda \sum_{(i,j) \in \mathcal{E}^{(k)}} E_{ij}(X_i^{(k)}, X_j^{(k)} : A_2), \quad (2.1)$$

where the first term is the data term indicating the cost of assigning a superpixel to foreground and background classes, while the second term is the smoothness term used for penalizing label disagreement between neighbours. Note that the $(:)$ part in these terms indicates that both these terms are functions of the learnt appearance model. From now on, to simplify notation, we write these terms as $E_i(X_i)$ and $E_{ij}(X_i, X_j)$, and the dependence on the appearance model \mathcal{A} and image (k) is implicit.

2.1.1 Data (Unary) Term

Our unary appearance model consists of a foreground and background Gaussian Mixture Model, *i.e.*, $A_1 = \{\text{GMM}_f, \text{GMM}_b\}$. Specifically, we extract colour features extracted from superpixels (as proposed by [14]). We use features from labelled sites in all images to fit foreground and background GMMs (where number of gaussians was automatically learnt by minimizing an MDL criteria [6]). We then use these learnt GMMs to compute the data terms for all sites, which is the negative log-likelihood of the features given the class model.

2.1.2 Smoothness (Pairwise) Term

The most commonly used smoothness term in energy minimization based segmentation methods [11, 12, 21] is the contrast sensitive Potts model:

$$E(X_i, X_j) = \mathbf{I}(X_i \neq X_j) \exp(-\beta d_{ij}), \quad (2.2)$$

where $I(\cdot)$ is an indicator function that is 1(0) if the input argument is true(false), d_{ij} is the distance between features at superpixels i and j and β is a scale parameter. Intuitively, this smoothness term tries to penalize label discontinuities among neighbouring sites but modulates the penalty via a contrast-sensitive term. Thus, if two adjacent superpixels are far apart in the feature space, there would be a smaller cost for assigning them different labels than if they were close. However, as various authors have noted, this contrast sensitive modulation forces the segmentation to follow strong edges in the image, which might not necessarily correspond to object boundaries. For example, [12] modulate the distance d_{ij} based on statistics of edge profile features learnt from a fully segmented training image.

In this work, we use a distance metric learning algorithm to *learn* these d_{ij} from user scribbles. The basic intuition is that when two features (which might be far apart in Euclidean distance) are both labelled as the same class by the user scribbles, we want the distance between them to be low. Similarly, when two features are labelled as different classes, we want the distance between them to be large, even if they happen to be close by in Euclidean space. Thus, this new distance metric captures the pairwise statistics of the data better than Euclidean distance. For example, if colours blue and white were both scribbled as foreground, then the new distance metric would learn a small distance between them, and thus, a blue-white edge in the image would be heavily penalized for label discontinuity, while the standard contrast sensitive model would not penalize this edge as much. The specific choice of this algorithms is not important, and any state-of-art technique may be used. We use the implementation of [5].

We update both $A_1 = \{\text{GMM}_f, \text{GMM}_b\}$ and $A_2 = \{d_{ij}\}$ every time the user provides a new scribble. Finally, we note that contrast-sensitive potts model leads to a submodular energy function. We use graph-cuts to efficiently compute the MAP labels for all images, using the publicly available implementations of [1] and [8,9,16].

2.1.3 Comparing Energy Functions

Our introduced energy functions (2.1) are different from those typically found in co-segmentation literature and we make the following observations. While previous works [13, 18, 19, 22] have formulated co-segmentation of image pairs with a single energy function, we assign to each image its own energy function. The reason we are able to do this is because we model the dependance between images implicitly via the common appearance model (\mathcal{A}), while previous works added an explicit histogram matching term to the common energy function. There are two distinct advantages of our approach. First, as several authors [13, 18, 19, 22] have pointed out, adding an explicit histogram matching term makes the energy function intractable. On the other hand, each one of our energy functions is submodular and can be solved with a single graph-cut. Second, this common energy function grows at least quadratically with the number of images in the group, making these approaches almost impossible to scale to dozens of images in a group. On the other hand, given

the appearance models, our collection of energy functions are completely independent. Thus the size of our problem only grows linearly in the number of images in the group, which is critical for interactive applications. In fact, each one of our energy functions may be optimized in parallel, making our approach amenable to distributed systems and multi-core architectures. Videos embedded on our project website [2] show our (single-core) implementation co-segmenting ~ 20 image in a matter of seconds.

To be fair, we should note that what allows us to set-up an efficiently solvable energy function is our incorporation of a user in the co-segmentation process, giving us partially labelled data (scribbles). While this user involvement is necessary because we work with globally related images, this involvement also means that the co-segmentation algorithm must be able to query/guide user scribbles, because users cannot be expected to examine all cutouts at each iteration. This is described next.

2.2 iCoseg: Guiding User Scribbles

In this section, we develop an intelligent recommendation algorithm to automatically seek user-scribbles and reduce the user effort. Given a set of initial scribbles from the user, we compute a recommendation map for each image in the group. The image (and region) with the highest recommendation score is presented to the user to receive more scribbles. Instead of committing to a single confusion measure as our recommendation score, which might be noisy, we use a number of “cues”. These cues are then combined to form a final recommendation map, as seen in [Fig. 2.1](#). The three categories of cues we use, and our approach to learning the weights of the combination are described next.

2.2.1 Uncertainty-based Cues

Node Uncertainty (NU): Our first cue is the one most commonly used in uncertainty sampling, *i.e.*, entropy of the node beliefs. Recall that each time scribbles are received, we fit $A_1 = \{\text{GMM}_f, \text{GMM}_b\}$ to the labelled superpixel features. Using this learnt A_1 , for each superpixel we normalize the foreground and background likelihoods to get a 2-class distribution and then compute the entropy of this distribution. The intuition behind this cue is that the more uniform the class distribution for a site, the more we would like to observe its label.

Edge Uncertainty (EU): The Query by Committee [23] algorithm is a fundamental work that forms the basis for many selective sampling works. The simple but elegant idea is to feed unlabelled data-points to a committee/set of classifiers and request

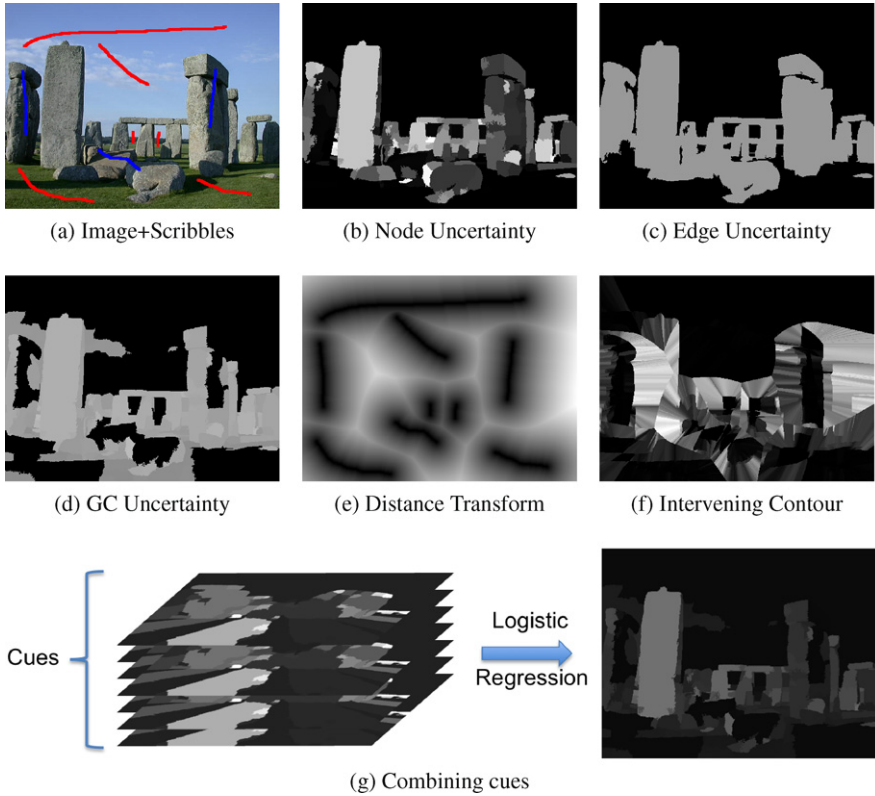


Fig. 2.1: Cues: (a) shows the image with provided scribbles; (b)-(f) show various cues; and (g) shows how these cues are combined to produce a final recommendation map.

label for the data-point with maximal disagreement among classifier outcomes. We use this intuition to define our next cue. For each superpixel, we use our learnt distances (recall: these are used to define the edge smoothness terms in our energy function) to find K ($=10$) nearest neighbours from the labelled superpixels. We treat the proportion of each class in the returned list as the probability of assigning that class to this site, and use the entropy of this distribution as our cue. The intuition behind this cue is that the more uniform this distribution, the more disagreement there is among the the returned neighbour labels, and the more we would like to observe the label of this site.

Graph-cut Uncertainty (GC): This cue tries to capture the confidence in the energy minimizing state returned by graph-cuts. For each site, we compute the increase in energy by flipping the optimal assignment at that site. The intuition behind this cue is that the smaller the energy difference by flipping the optimal assignment at a site,

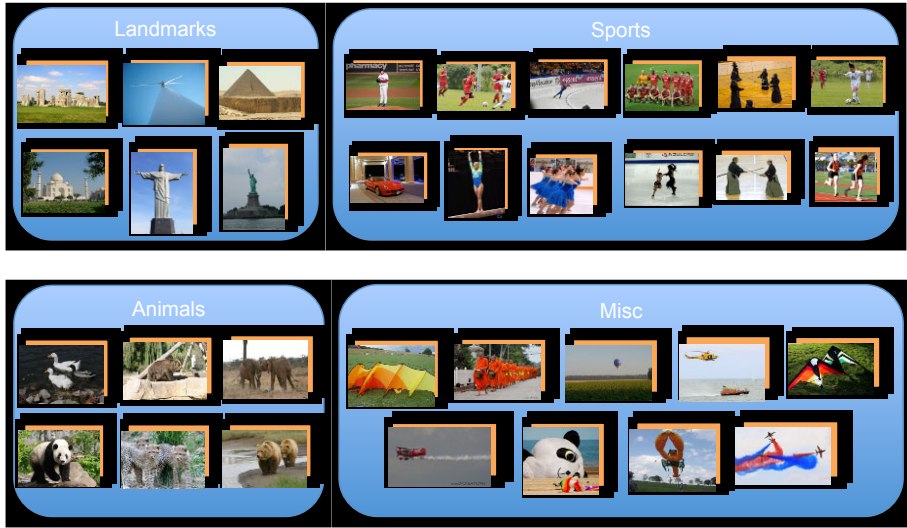


Fig. 2.2: CMU-Cornell iCoseg Dataset: Figure shows prototypical images from the co-segmentation groups in our dataset. Note: each image shown above corresponds to a *group* of images.

the more uncertain the system is of its label. We note that min-marginals proposed by [15] could also be used.

2.2.2 Scribble-based Cues

Distance Transform over Scribbles (DT): For this cue, we compute the distance of every pixel to the nearest scribble location. The intuition behind this (weak) cue is that we would like to explore regions in the image away from the current scribble because they hold potentially different features than sites closer to the current scribbles.

Intervening Contours over Scribbles (IC): This cue uses the idea of intervening contours [17]. The value of this cue at each pixel is the maximum edge magnitude in the straight line to the closest scribble. This results in low confusions as we move away from a scribble until a strong edge is observed, and then higher confusions on the other side of the edge. The motivation behind this cue is that edges in images typically denote contrast change, and by observing scribble labels on both sides of an edge, we can learn whether or not to respect such edges for future segmentations.

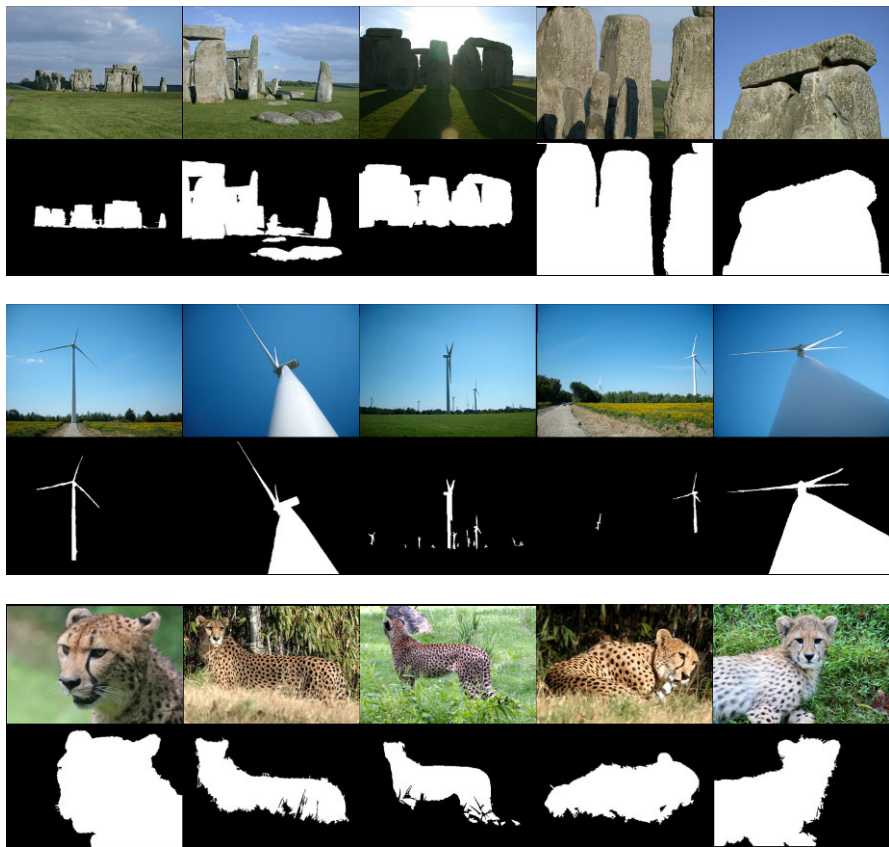


Fig. 2.3: CMU-Cornell iCoseg Dataset: Figure shows the Stonehenge, windmill and cheetah groups. We can see the large variation in illumination, scale, pose and appearance. Masks show the detail in our hand-annotated groundtruth.

2.2.3 Image-level Cues

The cues described so far, are local cues, that describe which region in an image should be scribbled on next. In addition to these, we also use some image-level cues (*i.e.*, uniform over an image), that help predict *which* image to scribble next, not where.

Segment size (SS): We observe that when very few scribbles are marked, energy minimization methods typically over-smooth and results in “whitewash” segmentations (entire image labelled as foreground or background). This cue incorporates a prior for balanced segmentations by assigning higher confusion scores to images with more skewed segmentations. We normalize the size of foreground and back-

ground regions to get class distributions for this image, and use the inverse of the entropy of this distribution as our cue.

Codeword Distribution over Images (CD): This image-level cue captures how diverse an image is, with the motivation being that scribbling on images containing more diversity among features would lead to better foreground/background models. To compute this cue, we cluster the features computed from all superpixels in the group to form a codebook, and the confusion score for each image is the entropy of the distribution over the codewords observed in the image. The intuition is that the more uniform the codeword distribution for an image the more diverse the appearances of different regions in the image.

2.2.4 Combined Recommendation Map

We now describe how we combine these various cues to produce a combined confusion map. Intuitively, the optimal combination scheme would be one that generates a recommendation map that assigns high values to regions that a user *would* scribble on, if they were to exhaustively examine all segmentations. Users typically scribble on regions that are incorrectly segmented. We cast the problem of learning the optimal set of weights for our cues as that of learning a mapping $\mathcal{F} : \phi_i \rightarrow \epsilon_i$, where ϕ_i is the 7-dimensional feature vector for superpixel i , corresponding to each of the 7 cues described above, and ϵ_i is the error indicator vector, which is 1 if the predicted segmentation at node i is incorrect, and 0 otherwise. We chose logistic regression as the form of this mapping. The ground-truth for training this logistic regression was generated by first scribbling on images³, co-segmenting based on these scribbles, and then using the mistakes (or the error-map) in these segmentations as the ground-truth. Our cue combination scheme is illustrated in Fig. 2.1.

2.3 The CMU-Cornell iCoseg Dataset

To evaluate our proposed approach and to establish a benchmark for future work, we introduce the largest co-segmentation dataset yet, the CMU-Cornell iCoseg Dataset. While previous works have experimented with a few pairs of images, our dataset contains 38 challenging groups with 643 total images (~ 17 images per group), with associated pixel-level ground truth. We built this dataset from the Flickr® online photo collection, and hand-labelled pixel-level segmentations in all images. We used the “Group” feature in Flickr, where users form groups around popular themes, to search for images from this theme. Our dataset consists of animals in the wild

³ More precisely, by generating random automatic scribbles on images. See Section 2.4.1 for details.

(elephants, pandas, *etc.*), popular landmarks (Taj Mahal, Stonehenge, *etc.*), sports teams (Baseball, Football, *etc.*) and other groups that contain a common theme or common foreground object. For some (though not all) of the groups, we restricted the images to come from the same photographer’s photo-stream, making this a more realistic scenario. Examples of these groups are shown in various figures in this paper, and Fig. 2.3 shows some prototypical images. We have made this dataset (and annotations) publicly available [2] to facilitate further work, and allow for easy comparisons.

Dataset Annotation: The ground-truth annotations for the dataset were manually generated by a single annotator using a labelling tool. The ground-truth was labelled on superpixels. However, the labelling tool allowed the annotator to interactively obtain finer / coarser superpixels as desired. A useful strategy used by our annotator was to use coarse superpixels while labelling simple scenes, and to use finer superpixels when labelling complicated or cluttered scenes. This allowed for very high quality ground truth, without the labeling task being prohibitively tedious.⁴

	# Groups	# Images	# Images/Group
[22]	7	16	2.29
[13]	23	46	2
CMU-Cornell iCoseg Dataset	38	643	16.92

Table 2.1: Dataset Statistics.

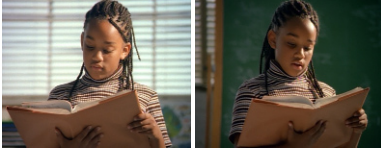
2.3.1 Dataset Statistics

We now analyze some statistics (size, appearance variation, scale variation) of our introduced dataset.

Size: Table 2.1 lists the number of groups, number of images and average number of images per group for our dataset. We note that this dataset is significantly larger than those used in previous works [13, 22]. Fig. 2.5a shows a histogram of number of images in groups.

Appearance: Recall our argument that when consumers take multiple photographs of the same event or object, the images are usually globally consistent. To quantify that images in our dataset do capture this property, we perform the following experiment. Fig. 2.4 shows a pair of images (“girl-pair”) from [22], and another

⁴ Since the superpixels were varied dynamically for each image, they were not the same as the ones used inside our co-segmentation algorithm (which only used a single setting of parameters for generating superpixels in all images).



(a) Girl-pair [22].

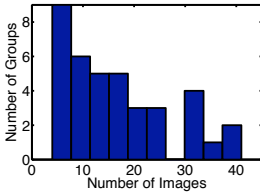


(b) Stonehenge-pair from CMU-Cornell iCoseg Dataset.

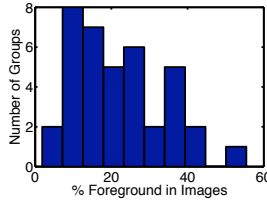
	KL-Divergence		
	Foreground	Background	Ratio
Girl-pair [22]	1.06	45.78	43.03
Stonehenge-pair	8.49	17.66	2.08

(c) Foreground and background similarity statistics

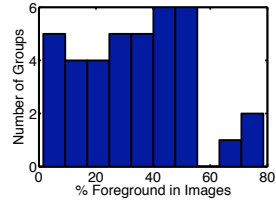
Fig. 2.4: Appearance Statistics: (a) shows a pair of images (“girl-pair”) from [22]; (b) shows a pair of images from the Stonehenge group in our CMU-Cornell iCoseg Dataset; (c) lists the KL-divergences between the two images for each pair. Images in our dataset are globally consistent, with comparable KL-divergence between foregrounds and backgrounds.



(a) Group Sizes.



(b) Hist. avg. foreground size



(c) largest - smallest foreground

Fig. 2.5: Dataset Statistics: (a) shows the histogram of the number of images in groups; (b) shows histogram of avg. foreground size in groups; (c) shows histogram of difference of largest and smallest foreground object within a group.

pair from our CMU-Cornell iCoseg Dataset (“Stonehenge-pair”). All the foreground pixels from the first pair of images were clustered into 64 color codewords using k-means clustering. A foreground histogram was built over this dictionary for each of the images, and the KL distance was computed between these normalized histograms. Similarly a color dictionary was built using the background pixels and the KL distance between the background distances was computed. This process was repeated for the second pair. We can see that for the “girl-pair” the appearance variation in foreground is considerably smaller than the variation in background. The “Stonehenge-pair”, on the other hand, shows a more comparable variation. This is not to say that the CMU-Cornell iCoseg dataset is inherently harder than previous datasets. Our intent is to point out that our dataset contains images where the pre-



Fig. 2.6: Scale Change: (a) shows the largest scale change, and (b) shows the smallest scale change in our dataset.

vious works on co-segmentation would fail, because they have been designed for a different scenario. We hope the introduction of our dataset motivates further research in the problem we consider in this paper.

Scale: To quantify the amount of scale change in our dataset, we show the histogram of average foreground size in groups in Fig. 2.5b. We can see that some groups contain very small foreground objects (on avg. $\leq 5\%$ of the image) while some groups contain very large foreground objects (on avg. $\geq 40\%$ of image). In addition, the histogram (Fig. 2.5c) of difference between largest and smallest foreground object in a group shows that even within a group there is significant scale change. Fig. 2.6 shows images with the largest and smallest scale change in our dataset. We can imagine that the large scale changes can be quite challenging for co-segmentation algorithms. In the case of our algorithm however, scale changes should not matter too much.

2.4 Experiments

For experimental evaluation, we performed machine experiments (Section 2.4.1) by generating synthetic scribbles, and also performed user-study (Section 2.4.2). In all experiments in this paper, we quantify the accuracy of an image segmentation as the percentage of pixels whose labels are correctly predicted. Co-segmentation accuracy for a group is the average segmentation accuracy over all images in this group.

2.4.1 Machine Experiments

To conduct a thorough set of experiments and evaluate various design choices, it is important to be able to perform multiple iterations without explicitly polling a human for scribbles. Thus, we develop a mechanism to generate automatic scribbles, that mimic human scribbles. We model the scribbles as (smooth) random walks that do not cross foreground-background boundaries. Our scribble generation technique consists of sampling a starting point in the image uniformly at random. A direction

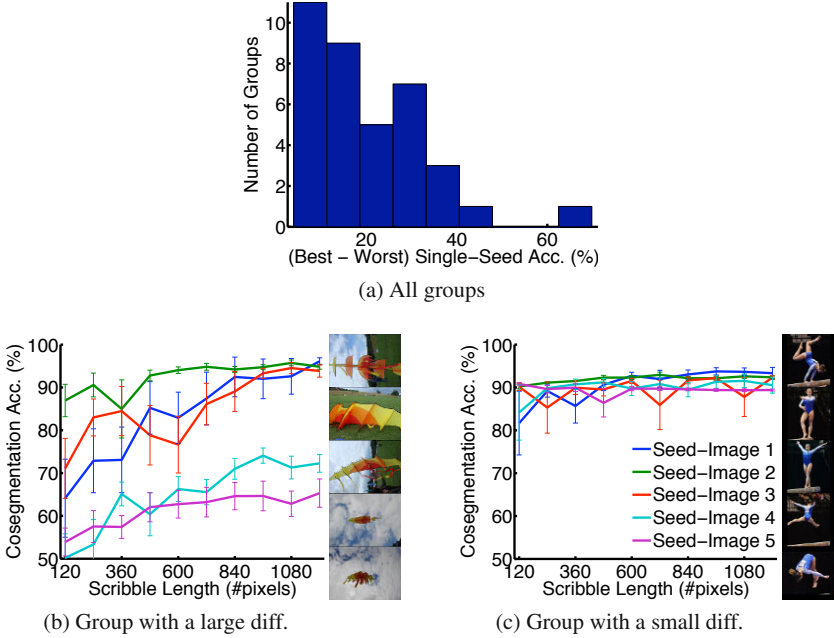


Fig. 2.7: Diversity within a group: (a) histogram of the difference in accuracy between the best and worst seed-images for all the groups in our dataset. A large difference indicates diversity in appearance. Some groups such as (b) “Kite” have images with varied appearances (two images providing worst segmentation accuracies don’t contain any grass), while other groups such as (c) “Gymnast” are more homogenous. Best viewed in colour.

angle is then randomly sampled such that it is highly correlated with the previous direction sample (for smoothness) for the scribble,⁵ and a fixed-size ($=30$ pixels) step is taken along this direction to extend the scribble (as long as it does not cross object boundaries, as indicated by the groundtruth segmentation of the image). To mimic user-scribbles given a recommendation map, the initial as well as subsequent points on the scribble are picked by considering the recommendation map to be a distribution. Using synthetic scribbles allows us to control the length of scribbles and observe the behavior of the algorithm with increasing information. Example synthetic scribbles are shown in Fig. 2.9. For all experiments in this paper, the length of each individual scribble was capped at 120 pixels.

⁵ For the first two sampled points, there is no previous direction and this direction is sampled uniformly at random.

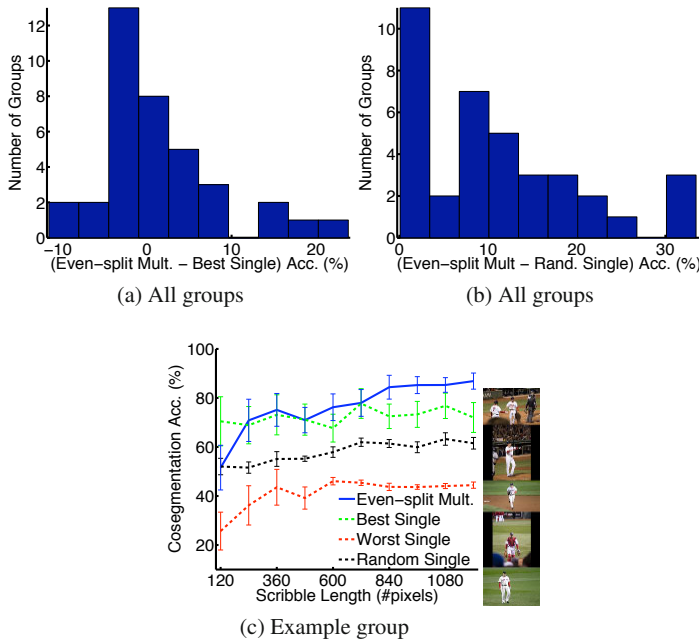


Fig. 2.8: Sparse vs. Dense Scribbles: (a) and (b) show the histograms of accuracy gains of even-split over the best (single) seed-image and a random (single) seed-image, respectively; (c) compares the accuracies achieved by dense scribbles on a single image, and sparse scribbles on multiple images for the group shown. Again, images are ordered in decreasing accuracy. The second histogram (b) shows providing scribbles on multiple images is a better strategy than committing to a single one.

2.4.1.1 Baseline 1: Scribbles restricted to a Single Image

To establish the simplest baseline, we ask the following question: “how well would interactive co-segmentation work if we were restricted to scribbling on a single image?” If a group consisted of successive frames from a video sequence, the choice of this chosen image (seed-image) would not matter much. The higher the diversity in the images among a group, the more variation we would observe in the group segmentation accuracies achieved by various seed-images, because not all seed-images would provide useful statistics for the group as a whole. We use the synthetically generated scribbles (described above) to test this.

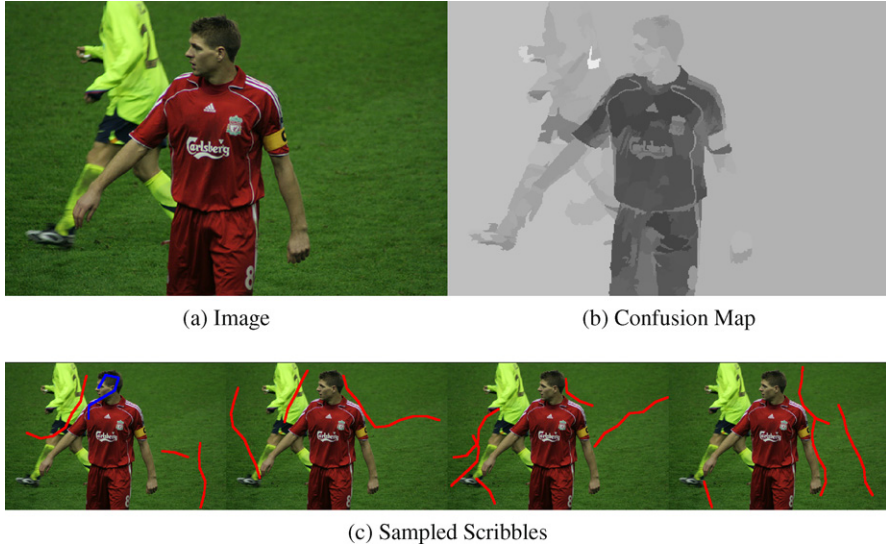


Fig. 2.9: Example simulated scribbles: Note that these scribbles never cross foreground boundaries (red player).

In Fig. 2.7a we show the histogram of the difference in accuracy between the best and worst seed-images for all the groups in our dataset.⁶ We can see that the histogram has a heavy tail, with 28 (of 37) groups having greater than 10% difference (and one as high as 69%), indicating that most groups have a lot of variation. Figs. 2.7b and 2.7c show the co-segmentation accuracies (Y-axis) for the “kite” and “gymnast” group respectively, as a function of the total length of scribbles (X-axis) on different seed-images in these groups (shown next to them). Images are ordered in decreasing accuracy. The “Kite” groups is an example of the groups with a lot of variation while the images in the “Gymnast” group have similar performances. Notice that in the “Kite” group, the two “bad” images do not contain any grass, and thus irrespective of the length of the scribbles on them, the algorithm will not know whether grass is considered foreground or background, and thus the group co-segmentation accuracy does not rise above 75%.

2.4.1.2 Aside: Automatic Seed Image Selection:

Before we describe the second baseline that we compare our interactive co-segmentation approach to, we take a small detour related to the previous baseline. Clearly, if we

⁶ In order to keep statistics comparable across groups, we select a random subset of 5 images from all groups in our dataset. One of our groups consisted of 4 images only, so all our results are reported on 37 groups.

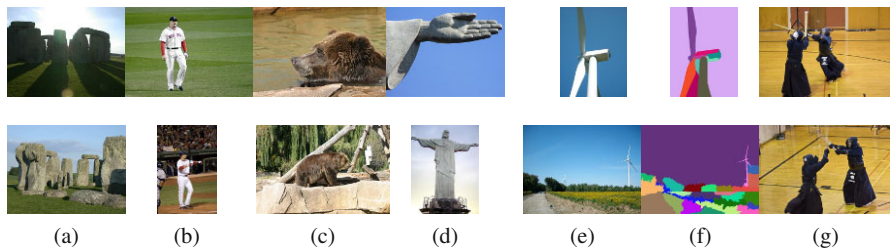


Fig. 2.10: For columns (a-e), top and bottom rows show best and worst images from example groups, which motivate our choice of features (a) Illumination histogram (b) HSV colour histogram (c) Gradient histogram (d) Gist and (e) Segmentation histogram; (f) Segmentations of images shown in (e); (g) Two images from the group are very similar with similar segmentation accuracies making the ranking a slightly misleading metric.

were restricted to scribble on only one image in a group, the choice of the seed image is crucial. This naturally leads to the question “Given a group of images to be co-segmented by scribbling on one image, can we automatically select this seed image in an intelligent way?”.

We pose this seed image selection problem as a classification task [4]. We extract the following features to describe an image.

- **Illumination histogram:** One of the observations we made was that the presence of strong shadows across the image (Fig. 2.10(a)) often results in poor co-segmentation accuracies. To capture this intuition, we compute a 50-dim histogram of the gray scale image.
- **Hue, Saturation and Value entropy:** The variety of colours in an image typically corresponds to the amount of useful information in the image, as seen in Fig. 2.10(b). We quantify this via a 3-dim vector holding the entropies of the hue, saturation and value marginal histograms. The more the number of colours present in an image, the higher the entropies in these distributions would be.
- **Gradient histogram:** The distribution of the strength of edges is a good indicator of how interesting the image is in terms of the existence of several regions/objects in the image. To represent this, we compute a 20-dim histogram of the edge magnitudes across the image. Fig. 2.10(c) shows that the worst image in the group has very few strong edges as compared to the best image which has more variety in its content.
- **Scene Gist:** The Gist features can help capture a holistic view of the overall scene layout (Fig. 2.10(d)). We extract the 1280-dim Gist features [20] which captures the response of the image to gabor filters of different orientations and scales, along with the spatial layout of these responses over the image.
- **Segmentation histogram:** Another indicator of the scene layout is the distribution of the sizes of segments in an image when run through an off-the-shelf segmen-

tation algorithm. For instance, as seen in Fig. 2.10(e,f), there is a stark contrast in the distribution of sizes of segments found in these images. We use mean-shift [10] for generating these segmentations.

We split our dataset into training groups and testing groups. We train a linear SVM using each of the n_f ($=5$) features described above individually to classify the best image in a group from the worst image. At training time we do not consider the remaining images in a group, because multiple images in groups can be visually similar leading to close cosegmentation accuracies, as seen in Fig. 2.10(g), and including them during training would make the classification problem artificially hard.

During testing, each image from the test group is passed through these n_f SVMs, and their output scores are recorded. Let the score corresponding to image x_i and SVM (feature) f_a be μ_i^a . The best images in training groups were labeled as the positive class, and thus we expect μ_i^a to be higher for better images.

We are ultimately interested in a ranking of the m images in the test group, and in order to do so, we compute a quality measure for each image, by comparing it to every other image in the group. Each image x_i is assigned a quality measure

$$Q(x_i) = \sum_{j=1}^m \sum_{a=1}^{n_f} [\![\mu_i^a - \mu_j^a]\!] \quad (2.3)$$

where $[\![t]\!]$ is the sign function, *i.e.* $+1$ if $t > 0$, and -1 otherwise.

This effectively captures how many times the image x_i got voted as being better than other images in the group, among all features. The m images in a group are ranked by this measure, and the top ranked image is chosen as the seed image.

For our experiments, we select m to be 5, and retain a random subset of 5 images from all groups. Since one of the 38 groups contained only 4 images, we work with the remaining 37 groups. We perform leave-one-out cross-validation on the groups. To understand the effectiveness of each of the individual features, we first report their corresponding image classification accuracies for identifying the best image from the worst. The results are shown in Fig. 2.11. It can be seen that all features hold some information to identify the best images from the worst ones (significantly outperforming chance, which would be 50%). It can be seen that the HSV entropies have the highest accuracy ($\sim 92\%$). This is understandable, especially since the segmentation algorithm uses colour features. All other features have similar accuracies ($\sim 76\%$).

To quantify the quality of the final ranking determined by our approach, we match our predicted ranks of images in the test groups, to the ground truth ranks (determined by sorting the average cosegmentation accuracies). We find that on average (across groups), we assign a rank of 2.14 to images that have a ground truth rank of 1. Moreover, the images that we select as rank 1, have, on average, a ground-truth rank of 2.11. In both cases, a random classifier would have an average rank of 3. Although this improvement in ranks may not seem significant, it should be noted that often groups contain more than one image that are “good” for scribbling and give similar segmentation accuracies, *e.g.* images shown in in Fig. 2.10(g).

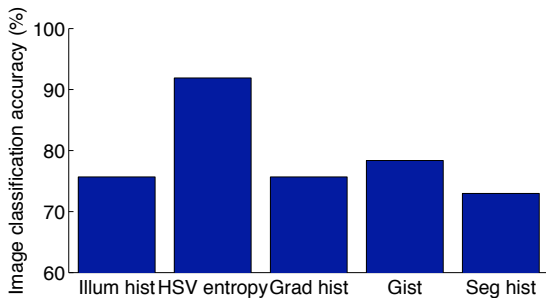


Fig. 2.11: The classification accuracies of each of our features in identifying best images in a group from the worst images.

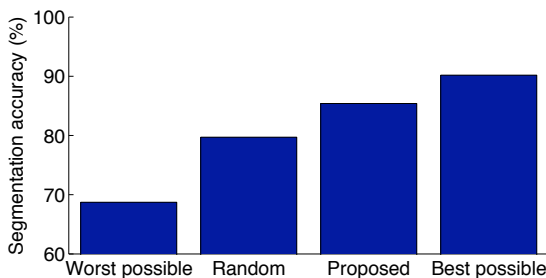


Fig. 2.12: The final cosegmentation accuracies.

The most relevant metric, for our application, is the gain in cosegmentation accuracies achieved by using our proposed seed image selection algorithm, as compared to picking an image from the group at random, which is the heuristic used by previous works [12, 24]. These results are shown in Fig. 2.12. We see that there is more than a 10% gap in the cosegmentation accuracies that can be achieved by scribbling on a randomly selected image (79.7%), and picking the best image in each group (90.2%). It should be noted that the best accuracy is the accuracy which would be achieved if an oracle were to label the best image in each group, and hence is the upper bound on what accuracy we can achieve. We can see that by scribbling on an image recommended by the seed image selection algorithm, we can fill more than half of this gap (at 85.4%).

Of course, the question of automatic seed image selection, while an interesting thought exercise, is relevant only if the user is restricted to scribbling on a single image in a group. As described thus far, our interactive co-segmentation system is more general, and recommends meaningful regions across all images in the group to the user to provide scribbles on. After this brief side analysis, we now focus again on the second baseline that we compare our approach to.

2.4.1.3 Baseline 2: Uniform Recommendation Maps

In the previous baseline we were restricted to scribbling on a single image. As we saw, the performance between the best and the worst seed-images could be significantly different. In this section, we consider user scribbles to be a limited resource and evaluate whether it is better to seek sparse scribbles on multiple images or dense scribbles on a single image. We follow a similar setup as in the last section, only now the scribbles are evenly split across all images in the group, which corresponds to uniform recommendation maps on all images. This way we can compare 1200-pixel scribbles (which would be dense) in a single image with five 240-pixel scribbles (which would be sparse).⁷ In practice, instead of making one long 1200-pixel scribble, we sample scribbles of length at most 120 pixels, and evenly split scribbles between foreground and background. As before, we perform 10 random runs. Fig. 2.8c shows the average cosegmentation accuracies in the group (Y-axis) for the worst (single) seed-image, the best (single) seed-image, a random (single) seed-image, and the accuracy achieved by evenly splitting scribbles across all images (called even-split) as a function of the total length of scribbles (X-axis). We can see that for the same length of scribbles, evenly splitting them across all images in the group and getting sparse scribbles performs better than dense scribbles on any image in this group. Fig. 2.8a and 2.8b show the histogram of accuracy gains of even-split over the best (single) seed-image and the random (single) seed-image experiments over all of the groups. The accuracies for Figs 2.8a, 2.8b were computed using scribbles of total length of 1200 pixels, *i.e.*, they correspond to the rightmost datapoint in Fig. 2.8c. We can see while even-split performs better than the best (single) seed-image for most groups, it is strictly better than a random (single) seed-image for *all* of the groups.

2.4.1.4 iCoseg

We first analyze the informativeness of each of our 7 cues. We start by generating a foreground and background scribble each of length at most 120 pixels on a random image in a group. We then compute each of our cues, and treat each individual cue as a recommendation map. We generate the next synthetic scribble (again of at most 120 pixels) as guided by this recommendation map, meaning that points are sampled by treating this recommendation map as a probability distribution (instead of sampling them randomly). We repeat this till we have scribbled about 1200 pixels across the group, and compute the average segmentation accuracy across the images of a group. We rank the 7 cues by this accuracy. Fig. 2.14 shows the mean ranks (across groups, average of 10 random runs) achieved by these cues. Out of our cues, the graph-cut cue (GC) performs the best, while both distance transform (DT) and intervening contour (IC) are the weakest. GC cue quantifies the uncertainty of the entire model (including node and edge potentials) and thus is expected to provide

⁷ This is one of the reasons for keeping a constant number of images per group. If each group had different images, even-split performance would no longer be comparable across groups.

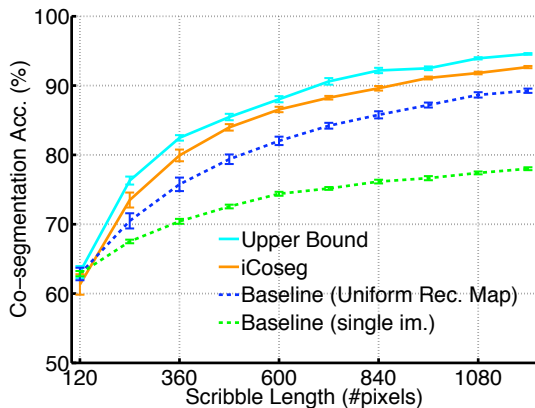


Fig. 2.13: Machine Experiments: Figure shows average co-segmentation accuracy as a function of the number of scribbles (each scribble is 120 pixels). iCoseg significantly outperforms baselines and is close to a natural upper-bound (see Section 2.4.1 for details).

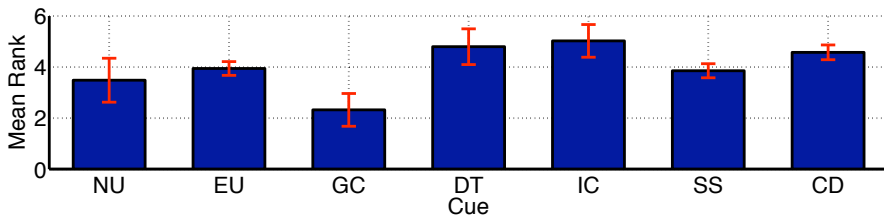


Fig. 2.14: Mean ranks achieved by individual cues (see Sec 2.4.1).

the best indication of where more information is required (from an active learning perspective). Thus, it is not surprising that this cue performs the best. DT and IC on the other hand completely ignore the learnt model, and only consider low-level cues like where (in x, y co-ordinates) we have scribbled in the image so far and the gradients in the image which often do not coincide with object boundaries. Thus, it is not surprising that they provide the least information to recommend meaningful regions to scribble further on.

We now evaluate iCoseg, our recommendation system, as a whole. The experimental set up is the same as that described above, except now we use the combined recommendation map to guide subsequent scribbles (and not individual cues). The cue combination weights are learnt from all groups except one that we test on (leave-one-out cross validation). We compare to two baselines described above. One is that of using a uniform recommendation map on all images in the group, which essentially means randomly scribbling on the images (respecting object boundaries of course). And the other (even weaker) baseline is that of selecting only one image

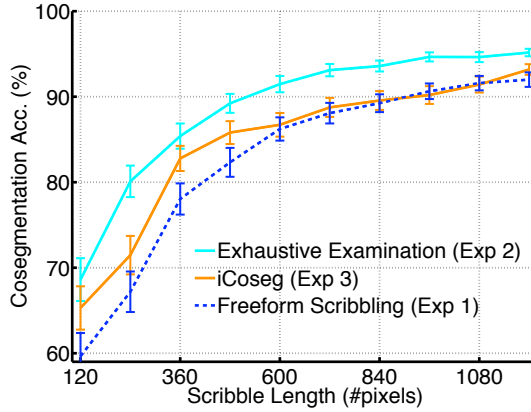


Fig. 2.15: User Study: Average performance of subjects in the three conducted experiments (see Section 2.4.2). iCoseg (Exp. 3) requires significantly less effort for users, and allows them to reach 80% co-seg accuracy with 75% the effort of Exp. 1.

(randomly) in a group to scribble on (with a uniform recommendation map on this image).

Fig. 2.13 shows the performance of our combined recommendation map (iCoseg) with increasing scribble length, as compared to the baselines. We see that our proposed recommendation scheme does in fact provide meaningful guidance for regions to be scribbled on next (as compared to the two baselines). A meaningful upper-bound would be the segmentation accuracy that could be achieved if an oracle told us where the segmentations were incorrect, and subsequent scribbles were provided only in these erroneous regions. As seen in Fig. 2.13, iCoseg performs very close to this upper bound, which means that users following our recommendations can achieve cutout performances comparable to those achieved by analyzing mistakes in all cutouts with significantly less effort *without* ever having to examine all cutouts explicitly.

2.4.2 User Study

In order to further test iCoseg, we developed a java-based user-interface for interactive co-segmentation.⁸ We conducted a user study to verify our hypothesis that our proposed approach can help *real* users produce good quality cutouts from a group of images, without needing to exhaustively examine mistakes in all images at each iteration. Our study involved 15 participants performing 3 experiments (each involving

⁸ We believe this interface may be useful to other researchers working on interactive applications and we have made it publicly available [2].

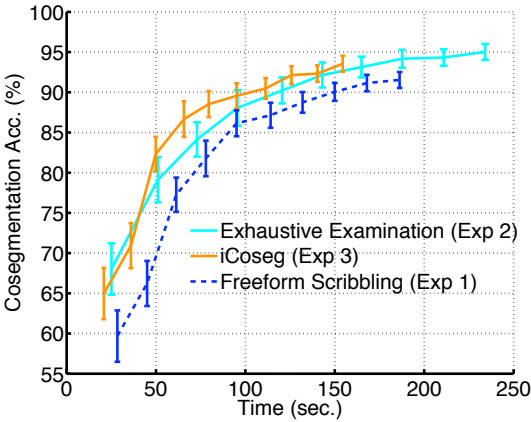


Fig. 2.16: User Study: Average performance of subjects in the three conducted experiments as a function of time (see Section 2.4.2). We can see that for a fixed amount of time, iCoseg (Exp. 3) typically achieves highest co-segmentation accuracy.

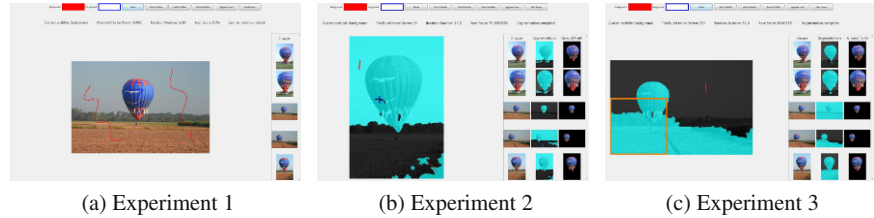


Fig. 2.17: User Study Screenshots: (a) Exp. 1: subjects were not shown cutouts and were free to scribble on any image/region while respecting the foreground/background boundaries; (b) Exp. 2: subjects exhaustively examine all segmentations and scribble on mistakes (cyan indicates foreground); (c) Exp. 3: users were instructed to scribble in the region recommended by iCoseg. Best viewed in colour.

	Mean no. of fg/bg scribbles per iter	Mean length of scribbles (px)	Mean % fg. in scribbles	Amount of uncertainty under the scribbles
Machine experiments	1.00 / 1.00	120	23%	-
User studies				
Exp. 1	1.04 / 1.05	106	46%	0.49 ± 0.04
Exp. 2	1.07 / 1.07	110	-	0.49 ± 0.04
Exp. 3	1.05 / 1.06	99	-	0.53 ± 0.04

Table 2.2: Comparison between user study and machine experiments.

5 groups of 5 related images). Fig. 2.17 shows screen-shots from the three experiments. The subjects were informed that the first experiment was to acclimatize them to the system. They could scribble anywhere on any image, as long as they used blue scribbles on foreground and red scribbles on background. The system computed cutouts based on their scribbles, but the subjects were never shown these cutouts. We call this experiment “freeform-scribbling”. In the second experiment, the subjects were shown the cutouts produced on all images in the group from their scribbles. Their goal was to achieve 95% co-segmentation accuracy in as few interactions as possible, and they could scribble on any image. We observed that a typical strategy used by subjects was to find the worst cutout at every iteration, and then add scribbles to correct it. In the third experiment, they had the same goal, but this time, while they were shown all cutouts, they were constrained to scribble within a window recommended by our algorithm, iCoseg. This window position was chosen by finding the location with the highest average recommendation value (in the combined recommendation map) in a neighbourhood of 201×201 pixels. The use of a window was merely to make the user-interface intuitive, and other choices could be explored. In all three experiments, users were restricted to use only 120 pixels of scribbles per iteration. Our UI displayed a counter that showed how many pixels they had left. Once their quota of pixels was over, they had no choice but to ask the system to co-segment using these scribbles, after which they were given a new quota of 120 pixels to scribble with. They did not have to use the entire quota before co-segmenting.

Fig. 2.15 shows the average segmentation accuracies achieved by the subjects in the three experiments (Y-axis) as a function of the length of their scribbles (X-axis). We can see that, as with the machine experiments, iCoseg helps the users perform better than freeform scribbling, in that the same segmentation accuracy (83%) can be achieved with about 75% the effort. In addition, the average time taken by the users for one iteration of scribbling reduced from 20.2 seconds (exhaustively examining all cutouts) to 14.2 seconds (iCoseg), an average saving of 60 seconds per group. Thus, our approach enables users to achieve cutout accuracies comparable to those achieved by analyzing mistakes in all cutouts, in significantly less time. This fact is further shown in Fig. 2.16 where the co-segmentation accuracy achieved (Y-axis) is plotted as a function of time taken (X-axis) for each of the three experiments, averaged across users and groups. We can see that our approach allows users to reach highest accuracies given the same time budget.

2.4.3 Comparing Machine Experiments and User Study

In order to understand how users scribbled in our user-study and study how well our automatic scribbles (machine experiments) emulate this, we analyze similarities between the user and synthetic scribbles. Table 2.2 compares some statistics.

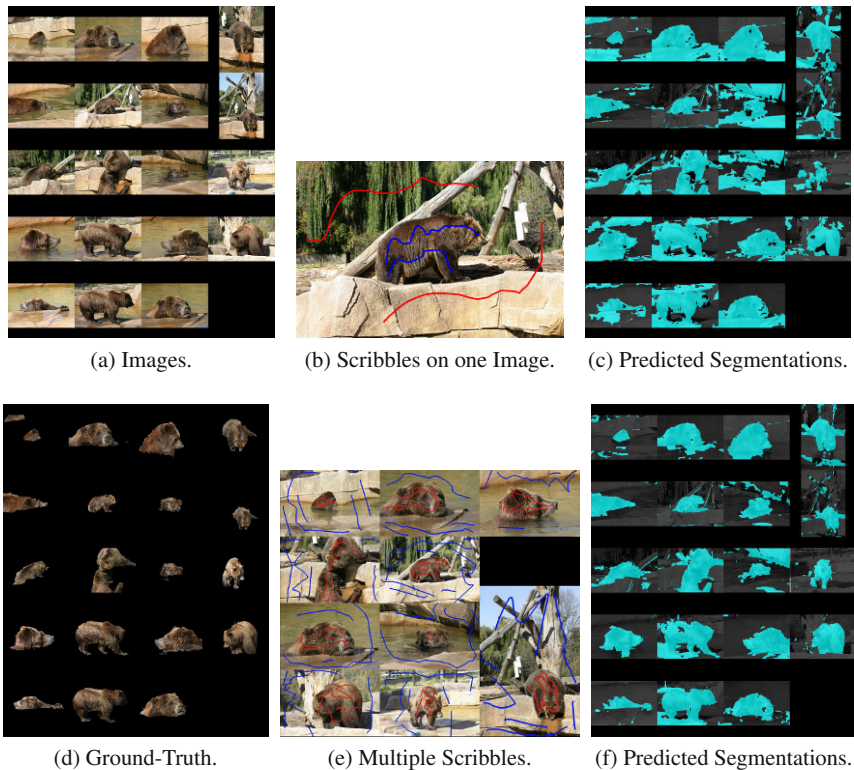


Fig. 2.18: Common Failure Case: (a) shows the group of images; (c) shows the segmentations achieved by scribbling on a single image, shown in (b). Cyan indicates foreground. (f) shows the segmentations achieved by scribbles on multiple images, shown in (e). When foreground and background have a lot of overlap in colour distributions, our interactive segmentation method faces difficulty in producing accurate segmentations (compare segmentations in (c),(f) with ground-truth in (d)). However, our algorithms allows for straightforward incorporation of more sophisticated features (*e.g.* colour-pallet of [12]), which should result in better performance.

Our automatic scribbles were generated for foreground and background with a fixed length of 120 pixels, and we see that they are comparable to the user scribbles in both the length and the average number of scribbles.

Interestingly, we also found that while our subjects were not given an explicit goal in Exp. 1 (freeform scribbling experiment) and were not shown the groundtruth, they were implicitly aware of the common foreground and their scribbles reflected that knowledge. The proportion of foreground pixels in all scribbles given by our subjects (for Exp. 1) was 46%, while the groups they viewed only contained 23% foreground. Clearly, they weren't scribbling uniformly randomly over an image, but were dividing their scribbles somewhat evenly over the foreground and back-

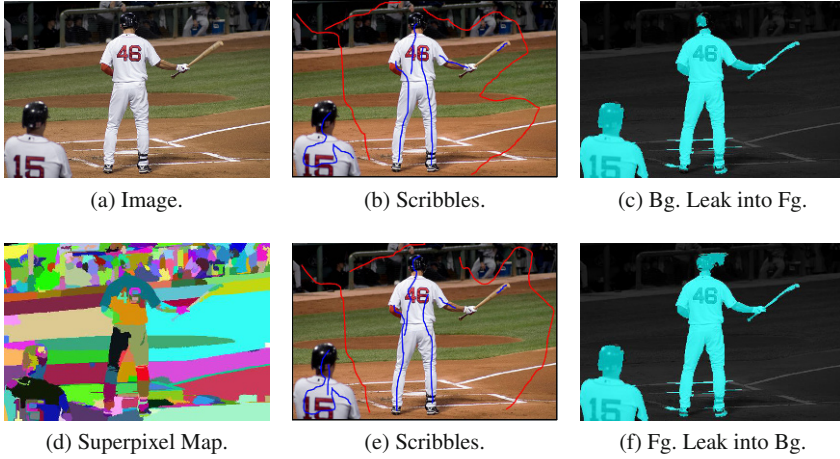


Fig. 2.19: Superpixel Leaks: We use a single parameter setting to generate superpixels for all images in our dataset, and thus some images show superpixel leaks across foreground objects. For example, in (d), the superpixel on the head of the baseball player leaks into the background. As a result, the segmentation also tends to either leak from foreground to background (f), or from background to foreground (c), depending on the choice of scribbles, (e) and (b) respectively.

ground. Thus, even though Fig. 2.15 seems to suggest that iCoseg has fallen to the performance level of the freeform-scribble baseline (Exp. 1), in reality, “freeform scribbling” has become a smart human-attention based algorithm. The truly random baseline can be seen in Fig. 2.13 where we force the random scribbles to be truly uniformly random, which our algorithm easily outperforms.

We also measured whether users were scribbling on confusing areas, as measured by our combined uncertainty map (which is normalized to be a spatial probability distribution, and thus between 0 and 1). We notice that the amount of uncertainty under user-scribbles is 0.49 ± 0.04 for both Exp. 1 (freeform scribbling) and Exp. 2 (exhaustive examination), again indicating that the users were implicitly aware about the common foreground and scribbled over incorrect segmentations which are typically regions with high uncertainty. We note that the uncertainty under user scribbles increased to 0.53 ± 0.04 for Exp. 3 (iCoseg), which is understandable, because the users were guided to scribble within the indicated regions of high uncertainty.

2.4.4 Limitations and Failure Cases

An assumption of our approach is that the foreground and background models are different enough in the chosen feature space (*i.e.* colour for our experiments) to

allow for reliable labelling of both classes. The interactive nature of our system makes the choice of features and appearance models seem less critical. However, they play an important role, and it is important to analyze the limitations and failure cases of our approach.

Non-discriminative Features: The most common failure case for our method results from the choice of colour features. Fig. 2.18 shows a difficult group to segment because the foreground colour distribution is very similar to the background colour distribution. Thus, even though the scribble guidance leads users to useful locations, the co-segmentation quality does not significantly improve despite multiple rounds of scribbles. Figs. 2.18b,c show the co-segmentations after scribbling on a single image, and Figs. 2.18e,f show the co-segmentations after scribbling on multiple images. We note that the choice of colour features is not inherent to the system, and more sophisticated features can be seamlessly incorporated. One choice of better features would be color-pattern features of [12] that capture the spatial distribution of colors in a neighborhood. These would provide more discriminative power (which should result in improved performance), as well as help overcome the local nature of features extracted at superpixels.

Superpixel leaks: We use superpixels as the labelling sites in our framework. This speeds up our implementation because the graph constructed on superpixels is significantly smaller than the grid-graph on pixels. However, because we use a single parameter setting to generate superpixels for all images in our dataset, some images show superpixel leaks across foreground objects. Fig. 2.19 shows an example image. Notice that some superpixels leak across object boundaries, *e.g.*, the one on the head of the baseball player. As a result of this superpixel leak, the segmentation also tends to either leak from foreground to background or from background to foreground, depending on the choice of scribbles. Having said this, our approach can be trivially extended to work with pixels, for applications that require highly accurate segmentations.

Single Background Model: It is conceivable that the use of multiple background models within a group could be beneficial. However, the more models we wish to build, the more scribbles we are likely to need from users for the models to be informative. In the extreme, in order to have one background model for every image in the group, we would need sufficient scribbles in all images in the group. This, to some extent, would defeat the purpose of having a co-segmentation system, where the goal is to leverage the fact that topically-related images share foreground and background statistics, and hence can be co-segmented, and need not be segmented individually. As seen in our examples, a large proportion of the images within a group do share similar backgrounds, a property that should be exploited when possible, but these properties are of course, application dependent.

In order to quantify the above intuition, we performed the following experiment. We performed co-segmentation with synthetic scribbles for the following three cases:

- **Multiple Models, Independent Segmentation (MMIS)**. In this case, the synthetic user scribbles on each of the five images in a group, and each image is independently segmented. Intuitively, this is equivalent to running a standard Grab-cut-like method on each image in the group, thus forcing the user to scribble on all images. The appearance models are not shared and the user is forced to scribble on all images to be segmented.
- **Multiple Models, Co-Segmentation (MMCS)**. In this case, the synthetic user again scribbles on a single image in the group, however now all images in the group are segmented by sharing the appearance model learnt from the single scribbled image. This is repeated by scribbling on all images in the group one at a time. As we have already observed in Section 2.4.1.3 and Fig. 2.8, we do not expect this combination to perform well.
- **Single Model, Co-Segmentation (SMCS)**. This is the case described in our machine experiments (Section 2.4.1 and Fig. 2.13), where the synthetic user scribbles on all images in the group. All images in the group are segmented by sharing the appearance model learnt from all the scribbled images.

	MMIS	MMCS	SMCS
Segmentation Accuracy	97.07 %	79.71 %	92.67 %

Table 2.3: Segmentation accuracies for various setups averaged across groups.

As we can see in Table 2.3, MMIS is the best thing to do, *i.e.* to scribble on every image and segment all images independently. However, this requires users to scribble on all images, which is not feasible for scenarios where the group contains many images. On the other hand, SMCS relieves the user from this constraint of scribbling on all images and as our machine experiments and user study show, the savings provided by iCoseg are crucial.

Now that we have presented a detailed description of iCoseg which allows users to interactively co-segment an object of interest from a group of images, in the following chapter, we discuss various exciting applications that are enabled by such a co-segmentation tool.

References

1. Bagon, S.: Matlab Wrapper for Graph Cut. <http://www.wisdom.weizmann.ac.il/~bagon> (2006).
2. Batra, D., Kowdle, A., Parikh, D., Tang, K., Chen, T.: Interactive Cosegmentation by Touch. <http://amp.ece.cornell.edu/projects/touch-coseg/> (2009).
3. Batra, D., Kowdle, A., Parikh, D., Luo, J., Chen, T.: Interactively Co-segmenting Topically Related Images with Intelligent Scribble Guidance. International Journal of Computer Vision (2011).

4. Batra, D., Parikh, D., Kowdle, A., Chen, T., Luo, J.: Seed Image Selection in Interactive Cosegmentation. International Conference on Image Processing (2009).
5. Batra, D., Sukthankar, R., Chen, T.: Semi-Supervised Clustering via Learnt Codeword Distances. British Machine Vision Conference (2008).
6. Bouman, C.A.: Cluster: An unsupervised algorithm for modeling Gaussian mixtures. Available from <http://www.ece.purdue.edu/~string-bouman>. (1997).
7. Boykov, Y.Y., Jolly, M.-P.: Interactive graph cuts for optimal boundary and region segmentation of objects in N-D images. International Conference on Computer Vision (2001).
8. Boykov, Y., Kolmogorov, V.: An Experimental Comparison of Min-Cut/Max-Flow Algorithms for Energy Minimization in Vision. IEEE Transactions on Pattern Analysis and Machine Intelligence (2004).
9. Boykov, Y., Veksler, O., Zabih, R.: Efficient Approximate Energy Minimization via Graph Cuts. IEEE Transactions on Pattern Analysis and Machine Intelligence (2001).
10. Comaniciu, D., Meer, P.: Mean shift: a robust approach toward feature space analysis. IEEE Transactions on Pattern Analysis and Machine Intelligence (2002).
11. Criminisi, A., Sharp, T., Blake, A.: GeoS: Geodesic Image Segmentation. European Conference on Computer Vision (2008).
12. Cui, J., Yang, Q., Wen, F., Wu, Q., Zhang, C., Gool, L.V., Tang, T.: Transductive Object Cutout. IEEE Conference on Computer Vision and Pattern Recognition (2008).
13. Hochbaum, D.S., Singh, V.: An efficient algorithm for co-segmentation. International Conference on Computer Vision (2009).
14. Hoiem, D., Efros, A.A., Hebert, M.: Geometric context from a single image. International Conference on Computer Vision (2005).
15. Kohli, P., Philip, T.H.S: Measuring uncertainty in graph cut solutions. Elsevier Journal on Computer Vision and Image Understanding (2008).
16. Kolmogorov, V., Zabih, R.: What Energy Functions can be Minimized via Graph Cuts?. IEEE Transactions on Pattern Analysis and Machine Intelligence (2004).
17. Leung, T., Malik, J.: Contour continuity in region based image segmentation. European Conference on Computer Vision (1998).
18. Mu, Y.D., Zhou, B.F.: Co-segmentation of Image Pairs with Quadratic Global Constraint in MRFs. Asian Conference on Computer Vision (2007).
19. Mukherjee, L., Singh, V., Dyer, C.R.: Half-integrality based algorithms for cosegmentation of images. IEEE Conference on Computer Vision and Pattern Recognition (2009).
20. Oliva, A., Torralba, A.: Modeling the Shape of the Scene: A Holistic Representation of the Spatial Envelope. International Journal of Computer Vision (2001).
21. Rother, C., Kolmogorov, V., Blake, A.: GrabCut: interactive foreground extraction using iterated graph cuts. ACM Transactions on Graphics (Proceedings of SIGGRAPH) (2004).
22. Rother, C., Minka, T., Blake, A., Kolmogorov, V.: Cosegmentation of Image Pairs by Histogram Matching - Incorporating a Global Constraint into MRFs (2006)
23. Seung, H.S., Oppor, M., Sompolinsky, H.: Query by committee. Computational Learning Theory (1992).
24. Schnitman, Y., Caspi, Y., Cohen O.D., Lischinski, D.: Inducing Semantic Segmentation from an Example. Asian Conference on Computer Vision (2006).

Interactive Co-segmentation of Objects in Image
Collections

Batra, D.; Kowdle, A.; Parikh, D.; Luo, J.; Chen, T.

2011, X, 46 p. 34 illus., Softcover

ISBN: 978-1-4614-1914-3