

# Research Statement

Saurabh Gupta

[sgupta@eecs.berkeley.edu](mailto:sgupta@eecs.berkeley.edu)

I am interested in computer vision, robotics and machine learning. My long term goal is to bring mobile robots into the real world. To be able to do so robots need to *perceive* the world and decide how to *act* on it to favorably change its state.

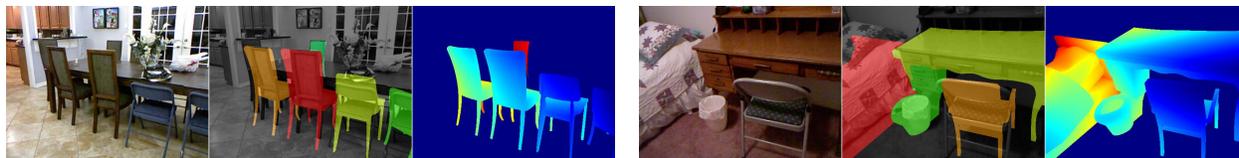
Building good actionable representations of the unstructured real world from raw sensory inputs is a fundamental question that needs to be addressed before mobile robots can be truly autonomous. My research has tackled this question and has led to the development of representations that can capture the structure of the world around us, can be efficiently learned using training data and can be derived from different sensing modalities. I have demonstrated the utility of these representations for *perceiving the world in 3D* and for *acting in the 3D world to move around*.

I started out by studying the problem of perception. Most work in computer vision focuses on 2D analysis of images. But our world is inherently 3D and a 2D analysis of the image is insufficient if we have to act in the world. This motivated me to pursue 3D scene understanding, a problem that has received much less attention. I first investigated this in context of range sensors (such as Microsoft Kinect). Here, I developed algorithms for analysis of range images to represent objects in the image with 3D models. This led to a detailed 3D understanding of real world scenes where we are able to make predictions for parts of the scene that are not even directly visible. I then extended this analysis to the case of RGB images, by developing algorithms to make inferences about 3D shape using just 2D images (Section 1).

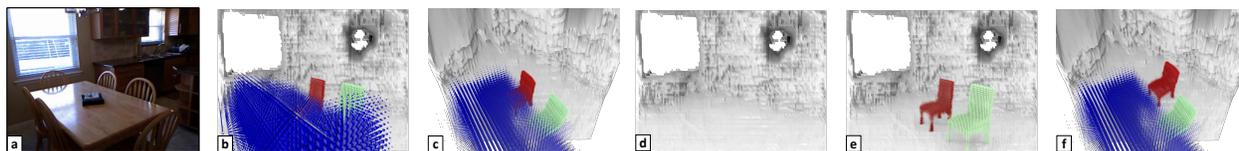
However, for robotics it is not enough to just perceive the world in 3D, it is equally important to be able to use this perception to take actions in the world. This led me to study the joint problem of perception and action, an area of research at the intersection of computer vision and robotics. Most existing work in this area falls in one of two paradigms: classical robotics and learning based robotics. Models in classical robotics use explicit hand-crafted geometric representations that capture the problem structure but ignore semantics and can't be learned from data. In contrast, models in learning based robotics can learn good representations in context of the end task, but ignore the structure present in the problem. Thus, results have mostly been confined to simple tasks often in 2D worlds, and extremely specialized behavior that does not generalize to new environments. This motivated me to formulate a new paradigm, that operationalizes insights about problem structure (from classical methods) into learning formalism (from learning-based robotics) through use of specialized end-to-end trainable policy architectures that jointly map and plan. Thus, our resulting policies not only benefit from learned task-driven representations, but also leverage the problem structure appropriately, leading to effective performance in novel environments (Section 2).

While working on these robotics problems, I realized that very often it is desirable to instrument robots with additional sensing capabilities (such as tactile sensors for grasping and manipulation or LiDARs sensors for self-driving cars). While on the one hand using custom sensors greatly simplifies the task, the dominant paradigm for learning representation relies on supervised training using large amounts of labeled data. This reliance on labeled data for learning representations severely limits the information that can be derived from such additional sensors. This led me to design techniques for transferring supervision from well-labeled modalities to other impoverished modalities that lack well-labeled large-scale datasets (Section 3).

I believe we are now very close to the point where robotic agents can be robustly and safely made to move around in unstructured uninstrumented real world environments. As we think about a future with such embodied robotic agents around us, a number of new research problems emerge. Rather than passively analyzing datasets collected by humans, machine learning algorithms will have the opportunity to collect their own data by moving around. It will require designing algorithms that *know what they don't know* and can design real-world behaviors to *teach themselves what they don't know*. This is just the beginning of an exciting era of, what I call, *artificial embodied cognition*, and I am excited to be a part of it.



**Figure 1: Detailed 3D Scene Understanding Output** as obtained using input from RGB-D sensors [9, 13]. We have detected objects of interest and replaced them in-place with similar CAD models. Two examples are shown here. Left picture shows the input color image, the middle picture shows the model masks rendered on the color image and the right picture shows the scene as composed by predicted models.



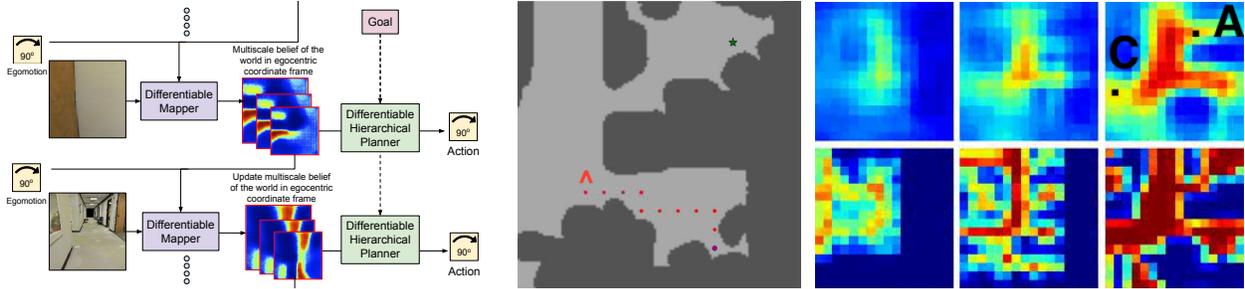
**Figure 2: Detailed 3D Scene Understanding Output** as obtained using input from RGB sensors [21]. As before we can detect objects and replace them in-place with voxel-based models. Additionally we can also predict the scene layout (floor, walls and ceiling). We demonstrate our understanding of the scene by removing objects and moving them around.

## 1 Detailed 3D Scene Understanding

Scene understanding is a central problem in computer vision and involves associating semantics with raw image observations. However, most computer vision algorithms only produce a crude understanding of the world and merely output a 2D bounding box around objects or at best infer their pixel support. In my work I have developed algorithms that go beyond such a 2D understanding of objects in an image to a full 3D understanding of objects in the world, including their complete 3D shape, location, pose and full unoccluded extent [9, 21]. This work has also unified three independent strands of research in computer vision: perceptual grouping (chunking raw pixels in an image into discrete entities), recognition (associating semantics with these groups of pixels), and 3D shape inference (inferring the 3D shape of these objects).

I have studied this problem in a variety of different contexts with challenging data captured from the real world. My early work used images from RGB-D sensors (such as Microsoft Kinect) where I designed algorithms for analyzing depth images. Not only is this setting directly relevant in context of robots that can be instrumented with additional sensors to boost performance, it also allowed us to address the scientific question about the role of observed 2.5D shape for the tasks of grouping and recognition. I first developed algorithms for perceptual grouping where I showed that using depth observations leads to large improvements for the specific tasks of contour detection and region proposal generation [8, 10, 13]. I was able to decrease the number of candidate regions required to cover objects in an image by an order of magnitude. I also showed how we can go beyond performance improvements on existing grouping tasks to more nuanced tasks such as amodal reasoning of scene surfaces, and classification of contours into different types (occlusion, albedo and convex and concave shape edges) [8, 10]. I then designed encodings for depth images for use in convolutional neural networks (CNN). Here, I found that providing additional geocentric context in the form of per-pixel estimate of the height above ground and the angle with respect to gravity led to boosts in performance for recognition tasks [13]. Numerous follow-up works from other researchers have adopted this geocentric encoding for contour detection [22] and semantic segmentation [6, 7] in RGB-D images, 3D object detection in outdoor settings [2], as well as inferring 3D shape attributes from RGB images [5]. Finally, I have also designed alignment algorithms that can infer and align full 3D CAD models to raw 2.5D observations in cluttered real world scenes, and have shown how information from grouping and recognition can together aid in this shape inference task, even without any annotations for ground truth shape or alignment [9]. Sample outputs are shown in Figure 1.

In some recent work [21] (under submission and in collaboration with other lab members), we have extended such detailed reasoning to scenarios where we only have RGB images. Additionally, the output is no longer restricted to being one of the previously seen CAD models. Insights from grouping and recognition motivated a factored representation where we model the scene as a composition of amodal scene surface (such as walls, floor and ceiling) and discrete objects (such as chairs, tables, beds) placed at different



**Figure 3: Visual Navigation [11]. Left:** Our specialized policy architecture that has mapping and planning components. The mapper processes information from the first person view to write into a spatial memory. The spatial memory serves as a ‘map’ that is used with a planner to output actions that will convey the robot to its desired goal location. **Right:** Our trained navigation policy can make predictions for parts of the environment that haven’t even been observed. We can predict free space inside the room even though we have only observed the doorway.

locations in the scene. We rely on synthetic data and train CNNs that output the location, pose, full 3D shape (parameterized by voxels) for objects and amodal depth of scene surfaces in the image. Some sample outputs are shown in Figure 2. Our factored representation is able to produce crisp output as opposed to past approaches that are fundamentally limited to only producing blurred output.

We can now predict detailed 3D representation of scenes that comprises of amodal surfaces and discrete objects along with their shape, pose and location from single RGB images. This is a large step forward from existing crude representations for scenes that comprise of collection of 2D object bounding boxes.

## 2 Visual Navigation

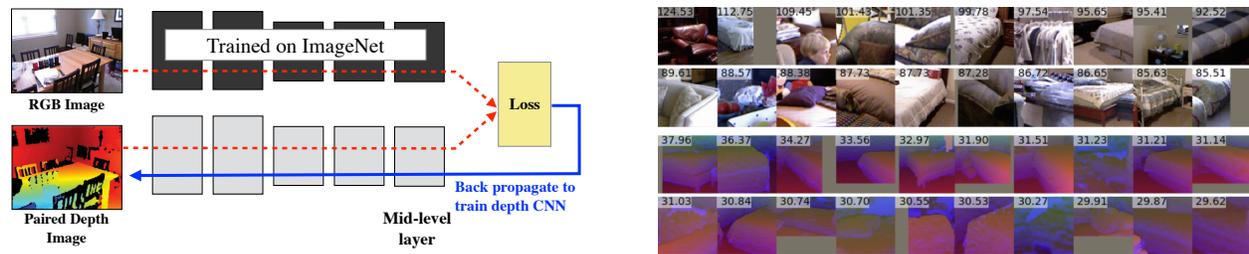
Robot navigation is an important problem in robotics and has been thoroughly studied over the last 40 years. The goal in robot navigation is to come up with control policies for robots that can allow them to efficiently move around the environment in a goal directed manner while avoiding collisions.

Classical approaches use purely geometric modeling of the world for inferring control policies. A purely geometric description of the world is neither a necessary nor a sufficient representation for navigation and fundamentally limits the capabilities and deployability of a robot. The precise reconstruction of a chandelier hanging from the ceiling may actually not be necessary for a ground robot to effectively move around. But more importantly an entirely geometric description ignores all semantics, limiting the capabilities of the robot. This manifests in a number of ways: knowing nothing till it is explicitly observed, and failure to exploit priors based on experience with similar past environments. This is in strong contrast with humans’ abilities to navigate. As humans, when we navigate through novel environments, we draw on our previous experience in similar conditions and are able to make meaningful speculations even about parts of the environment that have not been visited yet.

In my work [11] I have bridged this gap by formulating navigation as an end-to-end learning problem of choosing sequence of actions to optimally convey the robot to the goal location, given the sequence of first person observations (raw images) from on-board sensors. Such a learning problem is able to extract regularities of real world layouts that can then be used to make inferences based only on partial observations of new environments.

At the same time framing navigation as a learning problem is not enough by itself. As an agent moves around in the world, it gains experience of the particular environment it is in, allowing it to specialize its generic knowledge of spaces to the specifics of the current environment. Moreover, although the stream of observations being perceived by the agent is inherently sequential in nature, the underlying environment as well as its representation, even in simple rodents like rats, is inherently spatial (in the form of cognitive maps [20]). Additionally, it is difficult to specify how to use such spatial representations to determine the optimal action to reach the desired goal, specially when we want to allow them to naturally emerge through learning. This poses a problem for learning based methods where the typical unit of representation is unstructured feature vectors, sequence of which are treated using recurrent neural networks.

To overcome these challenges, I designed a trainable spatial memory module that accumulates information from the stream of observations to build a spatial representation of the world [11]. Moreover, instead of just



**Figure 4: Cross Modal Distillation for Supervision Transfer.** **Left:** An overview of the cross modal distillation technique for transferring supervision between modalities using paired data. **Right:** Aligned representations, corresponding neurons in the RGB and depth networks fire on semantically similar patches.

using the spatial representation as an unstructured feature vector in a classifier, I used it with a differentiable planner (based on the classical value iteration algorithm) to output optimal actions to reach the goal location. I also designed a version that builds and plans with these representations at different spatial scales leading to a very natural hierarchical decomposition of the problem. Figure 3 shows an overview of the designed policy architecture. Videos of these learned policies deployed on a robot are available on the project page.<sup>1</sup>

One of the challenges that comes about when maintaining and using maps is being able to localize oneself with respect to it (specially in contexts when there is noisy actuation). Classical approaches tackle this purely geometrically and rely on reasonably precise localization at all times. However, precise localization with respect to the map may actually not be necessary along all points on a trajectory, specially in context of the current observation and the rough heading direction (such as when going down a hallway, or when reaching for the door while exiting the room). Just as precise geometric mapping is unnecessary, precise localization at all points may be unnecessary. In some recent work [12] (currently under submission), I have tackled this issue by training a control policy that extracts out relevant information about the planned trajectory from past visual memories about the environment (and the map), and uses it in context of the current observations to move through environments under noisy actuation.

This design of navigation policies not only allows for learning patterns in real world, but also allows for specializing and improving the agent’s behavior as it becomes more experienced with the environment. But perhaps most importantly, our design showed how expressive representations that retain the intuition behind design of classical approaches can be learned with end-to-end techniques leading to a solution that benefits from the best of both worlds.

### 3 Multi-Modal Representations

The current dominant paradigm for obtaining representations involves supervised learning using large amounts of labeled data. This well for modalities for which it is possible to construct large-scale labeled datasets. For instance, this has worked very well in computer vision where representations are trained on the large scale ImageNet dataset and then adapted to the specific task of interest (say object detection or segmentation).

In robotics very often it is desirable to instrument the robot with additional sensing capabilities (such as LiDARs, SONARs, range, tactile sensors, microphones and infra-red cameras). While on the one hand using custom sensors can simplify the task significantly, the reliance on supervised learning for obtaining representations severely limits the information that can be derived from such additional sensors. The challenge is two folds, not only is there isn’t enough labeled data, there may not even be enough raw data to learn from. This poses an important question: how can representations be cheaply and efficiently learned for these additional modalities for which there aren’t large scale labeled datasets?

To answer this question, I designed a cross modal distillation technique for transferring supervision from well-labeled rich modalities (say RGB images) to other poor unlabeled sensing modalities (say depth images) using unlabeled paired data [14]. While it may be difficult to obtain a large amount of labeled data, it may be significantly easier to obtain paired data (by additionally adding a RGB camera in addition to the desired new sensor). Given such paired data we can set up a supervised learning problem where abstract semantic representations extracted from the rich modality (such as from high level CNN activations) are used as

<sup>1</sup><https://sites.google.com/view/cognitive-mapping-and-planning/>

targets for supervising learning of representations on paired poor modality (Figure 4 (left)). Such learning targets can in fact provide very dense supervision by exploiting spatial or temporal correspondence and implicit knowledge in CNN activations.

Thus, benefits of advances in training representations can be extended to a large number of other modalities simply by recording small amounts of paired data. Representations learned through cross modal distillation inherit all the usual properties of representations natively trained on a modality. Not only do we get a single good representation, we also get an entire hierarchy of representations of varying spatial resolution and complexity. Moreover, these representations keep improving as better RGB representations are developed. These representations can be further enhanced with training on target tasks with small amounts of labeled data. Such task specific training also makes these representations complimentary to the representation from the teacher modality and can in fact be used together to further boost performance. Finally, a by-product of this process is aligned representations (Figure 4 (right)) across different modalities that facilitates zero-shot transfer of models between modalities. This allowed us to transfer object detectors trained on large scale RGB object detection datasets to depth images, only using unlabeled paired data.

My technique for transferring supervision between modalities has been adopted by other researchers in computer vision and beyond. It has been used to obtain representations for different types of images (such as spherical images [19], low quality images [18]), other modalities such as sound [1], and also to transfer models between different languages [23].

## 4 Future Work

I believe we are very close to the point where robotic agents can robustly and safely be made to move around in unstructured uninstrumented real world environments. This ability by itself will be immensely useful and enable mobile home assistants, a natural next step after devices like Google Home and Amazon Echo. My research on how to represent the world (3D perception coupled with end-to-end learned action perception loops in multi-sensory setups) will play a key role in enabling this.

**Mobility in the Real World.** While these will be necessary components, by no means will they be sufficient. True large-scale mobility will require experience sharing between mobile agents of different form factors. In some ongoing work, we are learning control policies for collision avoidance that use the underlying dynamics of the robot. We have parameterized a model-based controller (iterative LQR) to follow a trajectory that goes to the desired goal location through a set of input way-points. We are training a policy to predict these way-points from first-person image observations (and desired goal) such that no collisions occur. This combination of parameterized model-based controllers, with policies for predicting the controller parameters, provides a setup for end-to-end learning that does not suffer from the exorbitant sample complexity of pure model-free methods. At the same time it also presents a setup for sharing experience among mobile agents of different form factors, necessary for enabling mobility at large scale.

Future research will also need to study navigation in dynamic environments with inanimate objects (that can be made to move) and animate other agents (that can move by their own will). Our factored object-centered representations are even more necessary for modeling such environments. At the same time memory representations that can track the state of the world over time will become more crucial as objects and agents undergo occlusions while they moving around. Beyond these technical challenges, a somewhat bigger question that will require interdisciplinary collaborations is how to design socially acceptable artificial mobile agents.

**Artificial Embodied Cognition.** As we think about a future with artificial robotic agents embodied in the real world, a number of new research problems emerge. Let us consider computer vision. Computer vision currently focuses on analyzing data captured by human agents and uploaded to the Internet. Humans determine what aspects of the real world are captured and consequently analyzed by computer vision algorithms. In contrast, as robots move around by themselves they will chose the view of the real world that is captured. The distribution of the data that computers will need to analyze will fundamentally change and this will require us to re-architect our approach to solving computer vision problems in a big way. Computer vision will change from being internet computer vision to *embodied computer vision*.

Current high performing computer vision systems rely on extensive manual annotations. As data gather-

ing becomes more and more automated with embodied mobile agents, it will no longer be feasible to annotate the data. Thus, unsupervised learning and cross-modal learning [14] will become increasingly important. Access to multiple modalities simultaneously plays a central role in bootstrapping learning in children [17]. It will be interesting to investigate if representations in computer vision can be similarly bootstrapped when given access to rich multi-sensory observation of the environment.

Embodied agents will also be able to generate their own supervision (such as by bootstrapping of existing models that work under specific viewpoints into models that can work from all viewpoints), leading to a paradigm shift in the way computer vision problems are thought of today. Moreover, mobile agents instrumented with multi-modal sensors can also collect large-scale datasets that are very hard to collect right now (*e.g.* about physical properties of the objects like coefficient of friction, elasticity, BRDFs).

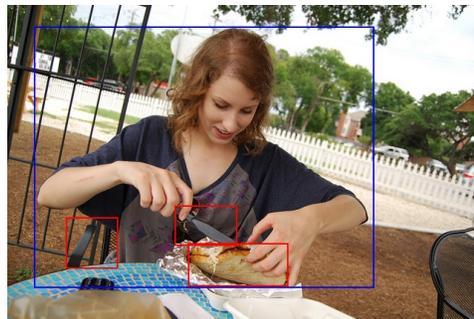
Mobile agents will also enable better understanding of human interactions with their environments. Observation of how humans change the state of the world around them is a rich information source for goal-directed behavior. This information can lead to better understanding of the world but is also a good starting point for imitation learning in context of robotics. In some early work, I studied human interactions with objects [16] and designed the task of visual semantic role labeling [15] (Figure 5). It will be great to revisit this work in context of home mobile agents that will naturally observe such data over time. 3D perception of the scene will be a natural aid as it can provide geometric context and occlusion reasoning for understanding interactions.

More broadly, artificial embodied cognition will require us to bring together insights from multiple different areas of artificial intelligence. My own work has benefited from knowledge of robotics, cognitive science [11], perception [13], machine learning [14], natural language understanding [3, 4, 15] and control theory (in ongoing work), leading to projects at the intersection of these different areas. I hope to continue to look at problems at these interesting intersections in my future research.

This is an extremely exciting time to be working on artificial intelligence. I think the best is yet to come, and I look forward to it.

## References

- [1] Y. Aytar, C. Vondrick, and A. Torralba. SoundNet: Learning sound representations from unlabeled video. In *NIPS*, 2016.
- [2] X. Chen, K. Kundu, Y. Zhu, H. Ma, S. Fidler, and R. Urtasun. 3D object proposals using stereo imagery for accurate object class detection. *PAMI*, 2017.
- [3] J. Devlin, H. Cheng, H. Fang, **S. Gupta**, L. Deng, X. He, G. Zweig, and M. Mitchell. Language models for image captioning: The quirks and what works. In *ACL*, 2015.
- [4] H. Fang\*, **S. Gupta\***, F. N. Iandola\*, R. Srivastava\*, L. Deng, P. Dollár, J. Gao, X. He, M. Mitchell, J. C. Platt, C. L. Zitnick, and G. Zweig. From captions to visual concepts and back. In *CVPR*, 2015.
- [5] D. F. Fouhey, A. Gupta, and A. Zisserman. 3D shape attributes. In *CVPR*, 2016.
- [6] A. Handa, V. Patraucean, V. Badrinarayanan, S. Stent, and R. Cipolla. Understanding real world indoor scenes with synthetic data. In *CVPR*, 2016.
- [7] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015.
- [8] **S. Gupta**, P. Arbeláez, R. Girshick, and J. Malik. Indoor scene understanding with RGB-D images: Bottom-up segmentation, object detection and semantic segmentation. *IJCV*, 2014.
- [9] **S. Gupta**, P. Arbeláez, R. Girshick, and J. Malik. Aligning 3D models to RGB-D images of cluttered scenes. In *CVPR*, 2015.
- [10] **S. Gupta**, P. Arbeláez, and J. Malik. Perceptual organization and recognition of indoor scenes from RGB-D images. In *CVPR*. 2013.
- [11] **S. Gupta**, J. Davidson, S. Levine, R. Sukthankar, and J. Malik. Cognitive mapping and planning for visual navigation. In *CVPR*, 2017.



**Figure 5: Visual Semantic Role Labeling:** Going beyond simple action classification to localizing the different semantic roles for different actions that the person is doing. Person is **cutting** bread with *knife* and **sitting** in a *chair*.

- [12] **S. Gupta**, D. Fouhey, S. Levine, and J. Malik. Unifying map and landmark based representations for visual navigation. *Under Review*, 2017. **Preprint**.
- [13] **S. Gupta**, R. Girshick, P. Arbeláez, and J. Malik. Learning rich features from RGB-D images for object detection and segmentation. In *ECCV*. 2014.
- [14] **S. Gupta**, J. Hoffman, and J. Malik. Cross modal distillation for supervision transfer. In *CVPR*, 2016.
- [15] **S. Gupta** and J. Malik. Visual semantic role labeling. *arXiv preprint arXiv:1505.04474*, 2015.
- [16] **S. Gupta\***, B. Hariharan\*, and J. Malik. Exploring person context and local scene context for object detection. *arXiv preprint arXiv:1511.08177*, 2015.
- [17] L. Smith and M. Gasser. The development of embodied cognition: Six lessons from babies. *Artificial life*, 2005.
- [18] J.-C. Su and S. Maji. Adapting models to signal degradation using distillation. In *BMVC*, 2017.
- [19] Y.-C. Su and K. Grauman. Learning spherical convolution for fast features from 360 imagery. In *NIPS*, 2017.
- [20] E. C. Tolman. Cognitive maps in rats and men. *Psychological review*, 1948.
- [21] S. Tulsiani, **S. Gupta**, D. Fouhey, A. Efros, and J. Malik. Factoring shape, pose, and layout from the 2D image of a 3D scene. *Under Review*, 2017. **Preprint**.
- [22] S. Xie and Z. Tu. Holistically-nested edge detection. In *ICCV*, 2015.
- [23] R. Xu and Y. Yang. Cross-lingual distillation for text classification. In *ACL*, 2017.