

Automatic localization and decoding of honeybee markers using deep convolutional neural networks

Benjamin Wild¹, Leon Sixt¹, Tim Landgraf^{1, *}

1 Dahlem Center of Machine Learning and Robotics, Freie Universität Berlin

* tim.landgraf@fu-berlin.de

Abstract

The honeybee is a fascinating model animal to investigate how collective behavior emerges from (inter-)actions of thousands of individuals. Bees may acquire unique memories throughout their lives. These experiences affect social interactions even over large time frames. Tracking and identifying all bees in the colony over their lifetimes therefore may likely shed light on the interplay of individual differences and colony behavior. This paper proposes a software pipeline based on two deep convolutional neural networks for the localization and decoding of custom binary markers that honeybees carry from their first to the last day in their life. We show that this approach outperforms similar systems proposed in recent literature. By opening this software for the public, we hope that the resulting datasets will help advancing the understanding of honeybee collective intelligence.

Introduction

Honeybees are a popular animal model in biology and have long served as inspiration in computer science. A honeybee colony can itself be seen as a distributed computing system. It manages numerous tasks in parallel with virtuosity. A bee colony adapts to significant environmental variation, it searches for and collects food, feeds its offspring, defends the hive against intruders and regulates temperature and humidity - all without central control. In the past, investigating the emergence of the abovementioned collective feats from individual behavior was mostly limited to synthetic approaches (Bonabeau et al. 1999, Beshers & Fewell 2001, Becher et al. 2014, Dornhaus et al. 2006). Analytical approaches through empirical studies have traditionally been limited in various dimensions. Numerous works have focused on the properties and effects of single behaviors, such as the famous waggle dance communication, that serves to direct and optimize the colony's foraging efforts (Grüter &

Farina 2009). The dance, however, is only one of the many communication channels over which nestmates may receive social information (Seeley 1995). Furthermore, due to the time-consuming nature of the biological experiment, virtually all studies were limited in the number of animals under observation, or its duration and sample rate (von Frisch 1965, Scheiner et al. 2013). v To make observational matters worse, bees are versatile learners. The life of a forager bee typically spans three to four weeks in which she may experience unique environmental features, such as rewards associated with olfactory cues, or locations that may bear perils. Previous work has shown that these memories can modulate a bee’s subsequent communication behavior. (Grüter et al. 2006, Balbuena et al. 2012, Grüter & Farina 2009, Grüter & Ratnieks 2011, Goyret & Farina 2005, De Marco et al. 2008, Richter & Waddington 1993, Nieh 2010). To fully understand how each of the colony’s individual members contribute to the collective, one arguably needs take into account each animals’ personal experience during their entire lifetime.

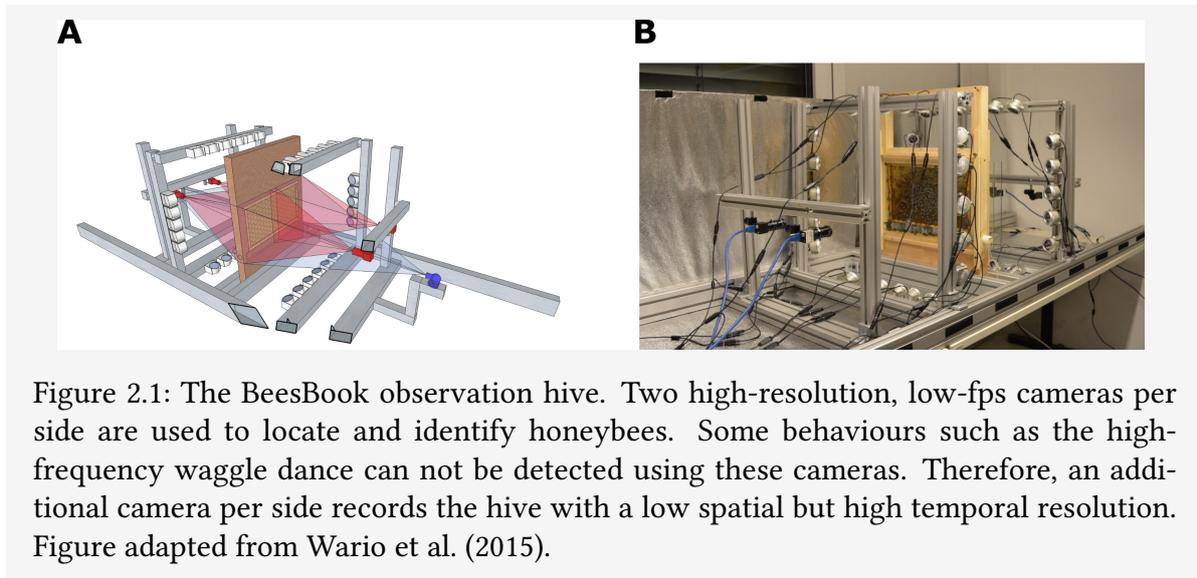
The behaviour of honeybees and other social insects was traditionally studied by manual observation and data collection (von Frisch 1965, Visscher & Seeley 1982, Seeley 1995). In the last decades, video recording technology was broadly adopted. On video, a larger number of individuals can be studied in detail (Beekman et al. 2004, De Marco et al. 2008, Landgraf et al. 2011, Wario et al. 2017). Effectively studying an entire bee colony over a long period of time, however, requires automation. Software for tracking unmarked animals is increasingly used in the behavioral sciences. Bees often leave the hive and even the young in-hive workers may be lost visually due to frequent occlusions. Because bees can not be identified based on their appearance alone, tracking systems for unmarked animals (Khan et al. 2004, Landgraf et al. 2007, Veeraraghavan et al. 2008) cannot be used for long-term observations.

In this paper, we describe a software pipeline for the automatic localization and identification of honeybees using custom binary identification markers. With this pipeline, we processed hundreds of terabytes of raw image data with high accuracy to create a unique dataset containing positions and orientations of all bees in a bee colony spanning several generations of workers.

State of the Art

Planar markers with binary codes have been shown to be feasible for tracking large groups of insects. A system, previously developed for ants (Mersch et al. 2013) was shown to successfully track 100 bees for two days (Blut et al. 2017). The markers used were originally described as fiducial markers in augmented reality systems (Fiala 2005) and rely on spatial derivatives

to detect the rectangular outline of a tag. A similar system using flat and rectangular markers for tracking larger insects was also proposed and might be adapted to honeybees (Crall et al. 2015). This system binarizes the image globally and searches for rectangular regions representing the corners of the marker. The previous BeesBook vision system (Wario et al. 2015) was tailored to specifically track all animals of small honeybee colonies over their entire lifetime. It uses a round and curved marker and searches for ellipse-shaped edge formations.



In Wario et al. 2015 we described a pipeline of conventional computer vision steps for detecting and decoding the markers. Although functional, this first prototype was computationally expensive and relied on a powerful supercomputer to process the large amounts of image data we recorded over three summer seasons. Furthermore, the decoding accuracy was drastically dependent on the image quality and respective parameters that had to be tuned for each of the four cameras used in the system.

In recent years, deep convolutional neural networks (DCNNs) advanced to the state of the art in many computer vision tasks (Krizhevsky et al. 2012, He et al. 2015, Razavian et al. 2014). Modern DCNN architectures such as VGG from the Visual Geometry Group (Simonyan & Zisserman 2014) and GoogLeNet (Szegedy et al. 2016) significantly outperform traditional image recognition techniques in essentially all common image classification benchmarks. Furthermore, deep convolutional neural networks are remarkably successful in object recognition and image segmentation tasks (Shelhamer et al. 2016, Girshick et al. 2013). Neural networks have been utilized to detect QR Codes (Chou et al. 2015) and to design application specific visual markers (Grinchuk et al. 2016).

we generated millions of image-labels-pairs clearing the way for using deep convolutional networks for decoding bee markers (see Figure 2.3).

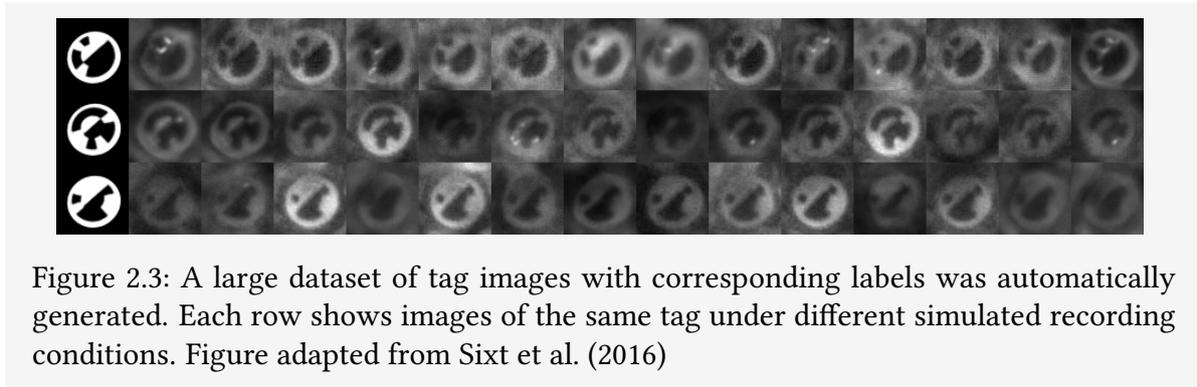


Figure 2.3: A large dataset of tag images with corresponding labels was automatically generated. Each row shows images of the same tag under different simulated recording conditions. Figure adapted from Sixt et al. (2016)

Methods

We identified the following requirements for the proposed machine vision solution. Since in the BeesBook project observations span several weeks, large volumes of image data have to be processed. Hence, the vision system needs to be fast and efficient. Honeybees populate the comb surface densely. The localization component therefore needs high spatial accuracy. Naturally, the detection and decoding accuracies should surpass the baseline described in Wario et al. (2015).

Execution time of convolutional neural networks is proportional to their size. Therefore, we propose using two separate models for marker localization and decoding, respectively. Tag localization does not require a high image resolution and may be a much easier task than decoding the marker. Therefore, a small fully convolutional net is used to localize the markers in a downsampled image. Only image regions containing bee markers are then processed in full resolution by a separate decoder network. This helps reducing processing time significantly.

Localization Positions of honeybee tags were manually annotated in a small subset of the raw BeesBook image data. To this end, the exact position of the center point of each tag was manually labeled in 179 images using a custom GUI. Because only the positions of the tags were required at this stage, this task took us only a few days.

Due to human error, we assume a normal distributed deviation of the center positions and propose using smoothly decaying labels rather than abstract coordinates as follows: Small

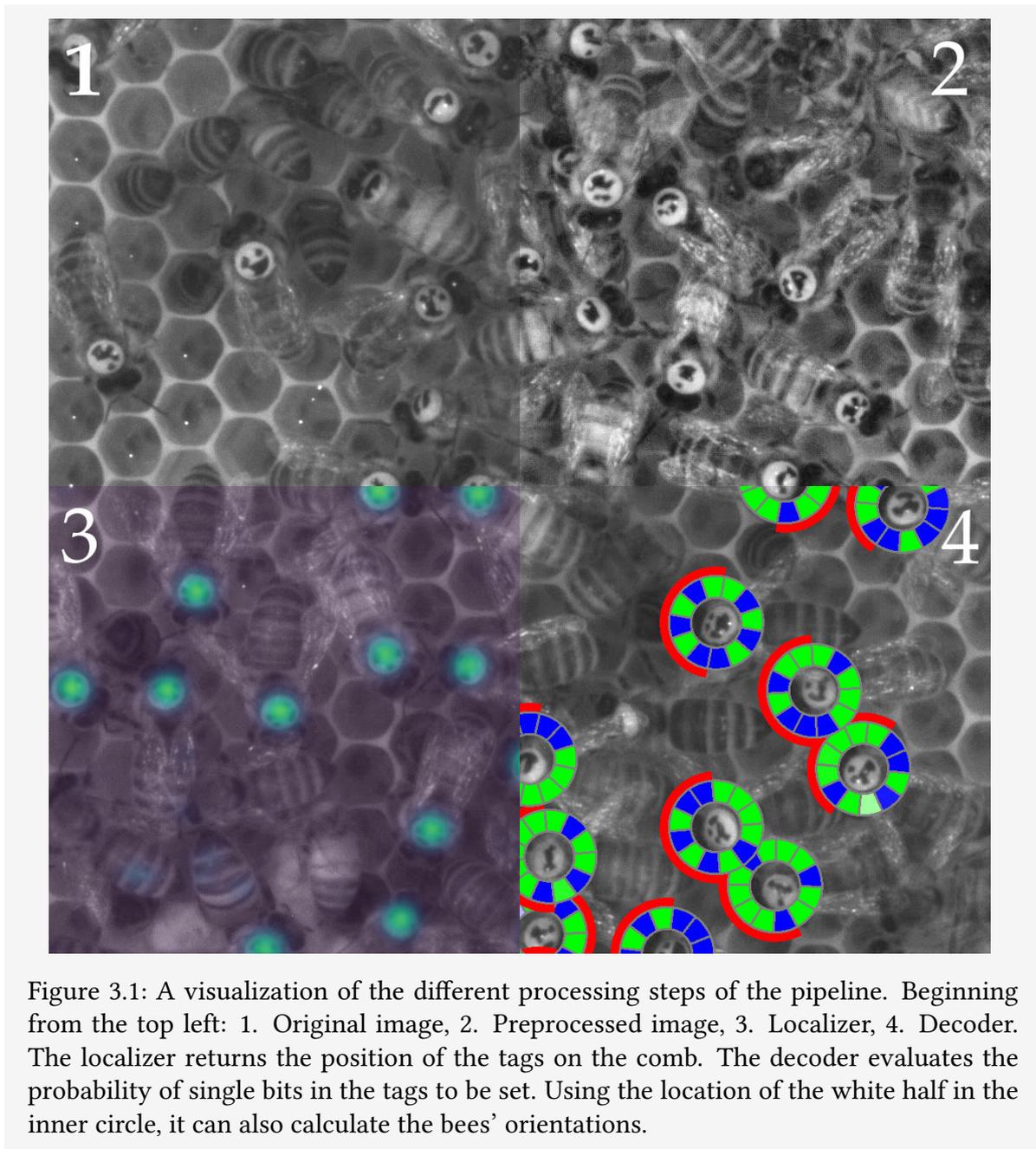
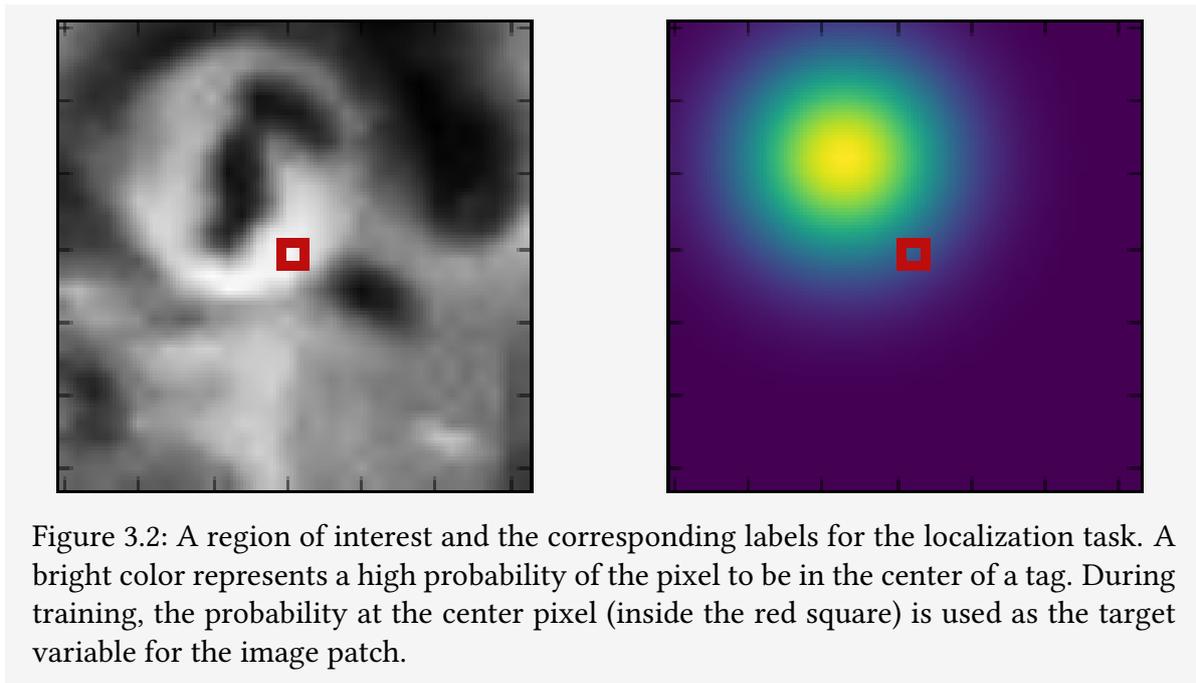


image patches (128 px width) were randomly sampled from the annotated ground truth dataset as inputs. We approximate the probability that a binary marker is exactly in the center of this image region with the density function of a bivariate normal distribution centered at the nearest true marker position with a fixed variance. The value of the density function at the center position of the image patch is used at the target variable in a regression setting (see Figure 3.2 for an example).



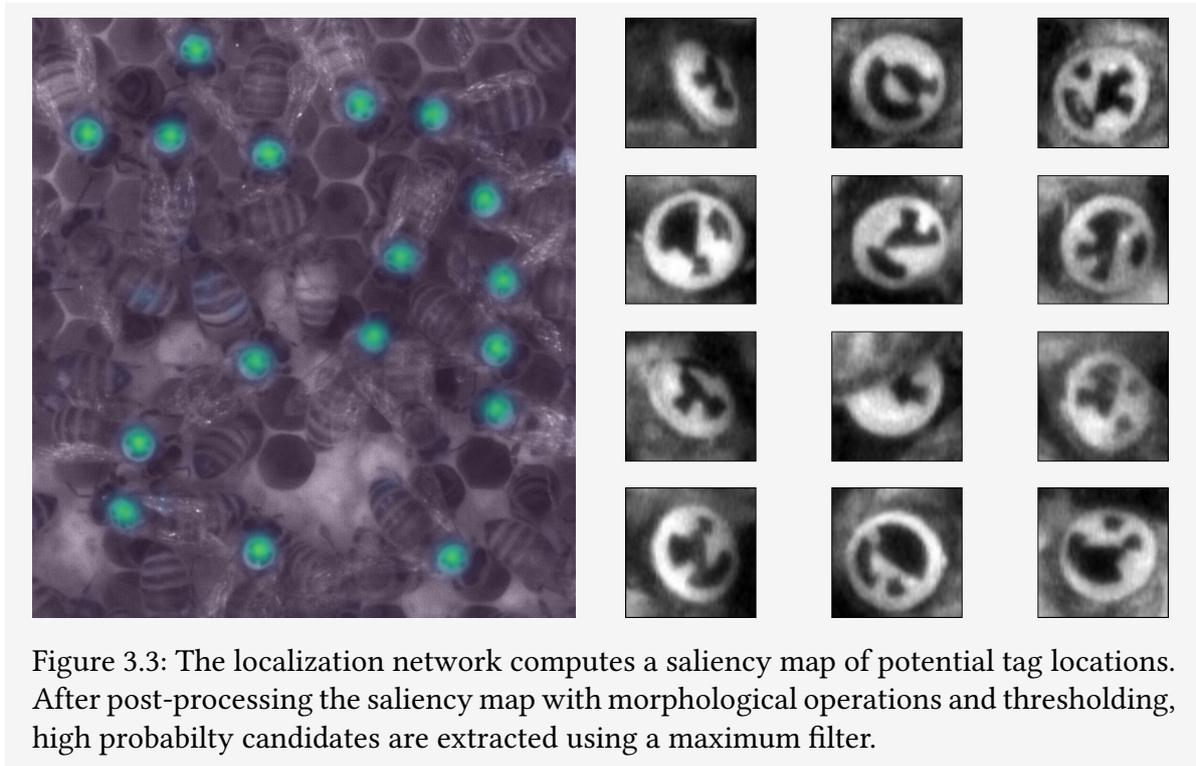
Using this method, a dataset of approximately 300.000 image patches (based on roughly 10.000 unique images of tag markers) was generated for the training of a marker localization model.

Images were preprocessed using the CLAHE (Pizer et al. 1987) algorithm to reduce the variance in brightness and contrast between different regions of the images, different cameras, and recording seasons. The image regions were then downsampled to a size of 32 px width using bilinear interpolation. Heavy data augmentation (such as random rotations and translations, added gaussian image noise, random elastic deformations, and random brightness and contrast perturbations) were used to prevent overfitting.

The localization model is a small fully convolutional neural network with three convolutional layers with a kernel size of 5 px (2 px stride) followed by a ReLU activation (Glorot et al. 2011). No padding is applied before the convolutional layers so that the model can be easily applied to full images during inference (and not only to the small regions of interest in the training dataset). Dropout (Srivastava et al. 2014) is applied following each convolutional layer to reduce overfitting. A final convolutional layer with a kernel size of 1 px (1 px stride) followed by a sigmoid activation computes the target probability for each output pixel.

The network is trained using stochastic gradient descent with momentum for a fixed number of 100 epochs.

During inference, we apply the network to compute a saliency map of the whole input image (after preprocessing and downsampling). Morphological operations combined with a maximum filter are then used to extract the local maxima. Maxima with a saliency below a fixed threshold are discarded. Image patches centered at these local maxima are extracted from the preprocessed images in full resolution and passed on to the next processing stages, e. g. the tag decoder or visualization.



Decoding While creating a labeled dataset for the localization task was manageable, it was unreasonably time-consuming for the decoding task. All members of our group used a custom GUI to visually match a three-dimensional grid onto all visible markers in a single camera recording. We obtained a labeled dataset of about 2000 marker instances in about one week time. This dataset served as evaluation set for the final system and not to train the convnet. Due to the complexity of the labels and the time-consuming nature of manual labeling, we developed a method to generate realistic images of markers that correctly display a given label. This work combines a typical GAN architecture (Goodfellow et al. 2014) with a 3D model of the bee markers. The GAN’s generator consists of a number of image augmentation functions that sequentially add image characteristics such as background, blur, noise and lighting

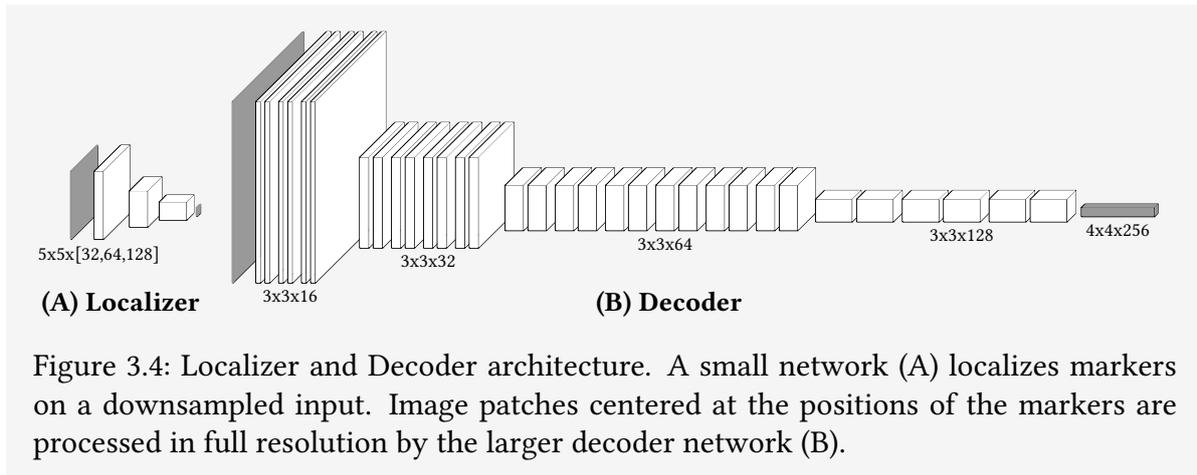


Figure 3.4: Localizer and Decoder architecture. A small network (A) localizes markers on a downsampled input. Image patches centered at the positions of the markers are processed in full resolution by the larger decoder network (B).

to synthetic images originating from the 3D model. The parameters of these augmentations are learned through adversarial training (Goodfellow et al. 2014) leveraging the previously mentioned database of marker images. The approach, called RenderGAN, allowed generating marker images with arbitrary combinations of spatial rotations and bit configurations. It was used to create a large dataset of 5 millions of labeled marker images for the training of a decoder model. See Sixt et al. (2016) for more details.

A large convolutional neural network based on the ResNet architecture (He et al. 2015) is used to decode the localized tags. Image patches containing markers with a size of 64 px are fed to an initial convolutional layer with a kernel size of 3 px (1 px stride) followed by a Batch Normalization (Ioffe & Szegedy 2015) layer and an ELU (Clevert et al. 2015) activation. The data is then processed following the 34-layer architecture described in He et al. (2015), but starting with only 16 filters in the first block to reduce the total numbers of parameters. The resulting representations are then individually processed by two fully connected layers with 256 filters each and followed by an ELU activation. The output from the first fully connected layer is finally used to compute the probabilities for each bit to be set using a sigmoid activation. The second fully connected layer is used to compute the spatial rotations of the marker (represented as vectors on the unit circle) and a full resolution offset for the exact center position of the tag. The bit probability outputs are optimized using a cross entropy loss while the mean squared error is used for the other outputs.

Potential decoding errors can still be corrected in later processing stages, e. g. during temporal tracking of the detections over time. To facilitate this, we store the raw probabilities for each bit instead of the thresholded predictions for each detection. Furthermore, we define a confidence measure for each prediction as:

$$c(b) = \frac{\prod_{i=0}^{12} 2 \cdot |0.5 - p(b_i)|}{12} \quad (1)$$

A high confidence measure $c(b)$ signifies that the decoder assigns a high probability to each of its predictions for the individual bits $p(b_i)$ which means that the decoding of the ID is likely to be correct. This measure can also be used to discard low confidence predictions in analyses where a high precision of the decoded IDs is more important than a high recall of detected bees.

Data storage While all data in the BeesBook project is ultimately stored in a database that can be used and extended easily by other team members, an intermediate data format was used to store the outputs of the pipeline due to the following reasons: Concurrently writing to a central database is difficult to achieve because of the peculiarities of a supercomputer’s job queueing system and creates a potential IO bottleneck (i. e. the workers can’t continue to process data because they have to wait for the database to store their previous results). Furthermore, the pipeline should work without any central coordination, i. e. it should be possible for any worker to process its local subset of the data without any communication to another process. Lastly, writing to the distributed filesystem on a supercomputer can be slow and therefore the size of the outputs should be as small as possible.

For these reasons, the results of the pipeline were serialized and saved using the Cap’n Proto library (Sandstorm Development Group 2013–2017). After processing, the results of the individual workers were collected and sorted by date in a simple file system hierarchy.

Results

We evaluated the performance of the two networks and compared them with our previous computer vision pipeline (Wario et al. 2015) and the bee tracking systems described in Blut et al. (2017) and Gernat et al. (2018). We also include the results of our approach when combined with temporal tracking as described in Boenisch et al. (2018).

Localizer The localizer convnet was trained to predict the saliency label as described in the previous section. During training, the model learns to minimize the binary cross-entropy between its outputs and the training labels. While this metric has proven to be very effective for training, it is less meaningful as a performance metric for detecting bees. We therefore calculate recall (the ratio of detected markers and the number of existing markers in the

same image) and precision (percentage of regions that correctly contained a marker). Both metrics depend on the threshold value used to discard local maxima with a low probability. A threshold of 0.6 was determined empirically and used in the evaluation and data processing. The new localizer outperforms the localization stage of the old computer vision pipeline in terms of runtime, recall, and precision without the need for hyperparameter tuning. The model can be applied without performance loss on images from all recording seasons and is able to recognize tags on images that are very different from the BeesBook observation hive recordings, for example tags on a white table photographed with a smartphone camera. Qualitative inspection of error cases suggests that the model sometimes misses tags which are very close to the image border (where the image sharpness is lower) and tags outside of the main surface area of the comb (where tags are usually farther away from the camera).

| | Recall | Precision |
|----------------------|---------------------|--------------|
| Wario et al. (2015) | $94 \pm 4\%$ | $88 \pm 4\%$ |
| Blut et al. (2017) | $90.8\% - 98.2\%$ * | – |
| Gernat et al. (2018) | $87 \pm 2\%$ | – |
| This work | 98.3% | 99.4% |

* Moving bees vs. resting bees

Figure 4.1: Comparison of recall and precision rates for localizing honeybee markers. The system we propose achieves best recall rates, i.e. our localizer misses the least amount of markers. In comparison to our previous system, we improved precision rates by 11%. In Blut et al. (2017) and Gernat et al. (2018), precision rates for localizing markers were not reported.

Decoder The decoder model not only predicts the values of the individual bits on the tag, but also all three spatial rotations. During training, the decoder learns to minimize the binary cross-entropy between the true values of the individual bits and its predictions and the mean squared error between the true orientations and its predictions.

We evaluate the mean hamming distance (the average number of bits decoded incorrectly) and the decoding accuracy (the number of tags that were decoded without any error).

Even though we do not use ECC, the decoder is able to decode 87.8% of the tags without errors. If we combine predictions from this DCNN with temporal tracking, the accuracy of the assigned IDs improves to 98.1% compared to 66% when using our computer vision pipeline (Boenisch et al. 2018).

| | MHD | Accuracy | Time per tag (ms) |
|------------------------------------|-----------|-------------|-------------------|
| Wario et al. (2015) | 1.08 | 66% | 177.54 |
| Gernat et al. (2018) | – | 98.58% | – |
| This work (w / wo tracking) | 0.42/0.08 | 87.8%/98.1% | 1.43/2.01 |

Figure 4.2: Comparison of marker decoding performance. Our systems do not use error correction schemes and may therefore exhibit flipped bits. The mean hamming distance (the expected number of flipped bits) was improved significantly. The proportion of correctly decoded markers (i. e. zero flipped bits) was improved by 21%. The system proposed in Gernat et al. (2018) surpasses our raw result by 11%. Applying a post-processing step to link corresponding detections to motion paths (Boenisch et al. 2018) allows us to match this decoding accuracy. Because of the improved performance, we can process our data in realtime on consumer hardware.

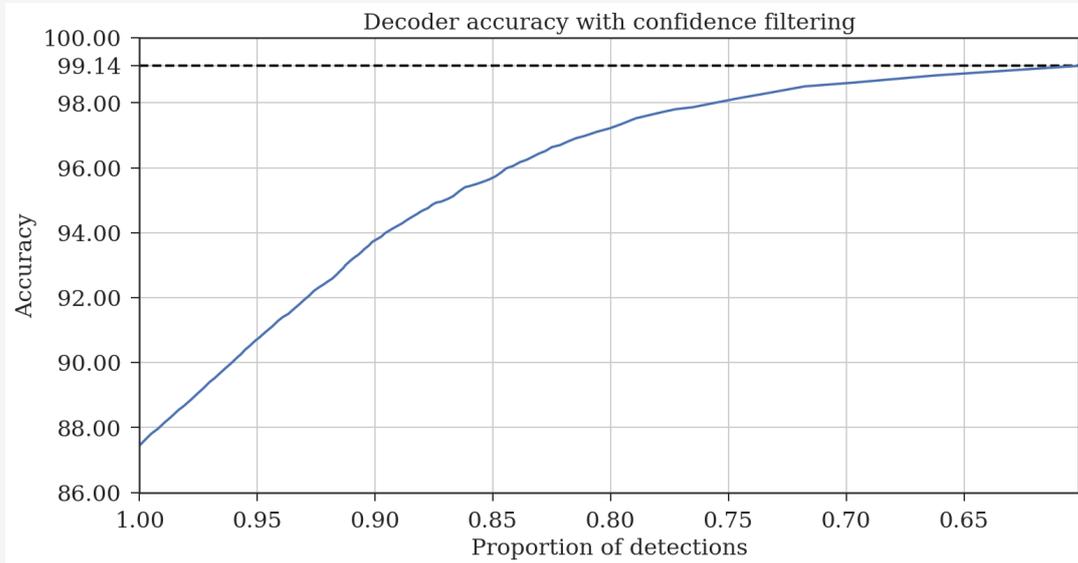


Figure 4.3: Confidence measure of decoder vs accuracy. Depending on the analysis, either a high recall or a high precision may be preferable. At maximum recall, the model can correctly decode approximately 88% of all detections. A very high decoding accuracy of 99.14% can be achieved without any temporal tracking by discarding the 40% of the detections with the lowest confidence score.

Resulting datasets All data of the recording seasons 2015 and 2016 was processed using the new deep learning pipeline on a Cray X30 supercomputer. The datasets are now available for further studies of the honeybee’s social behavior. Temporal tracking of the detections over time further improves the quality of the datasets (Boenisch et al. 2018).

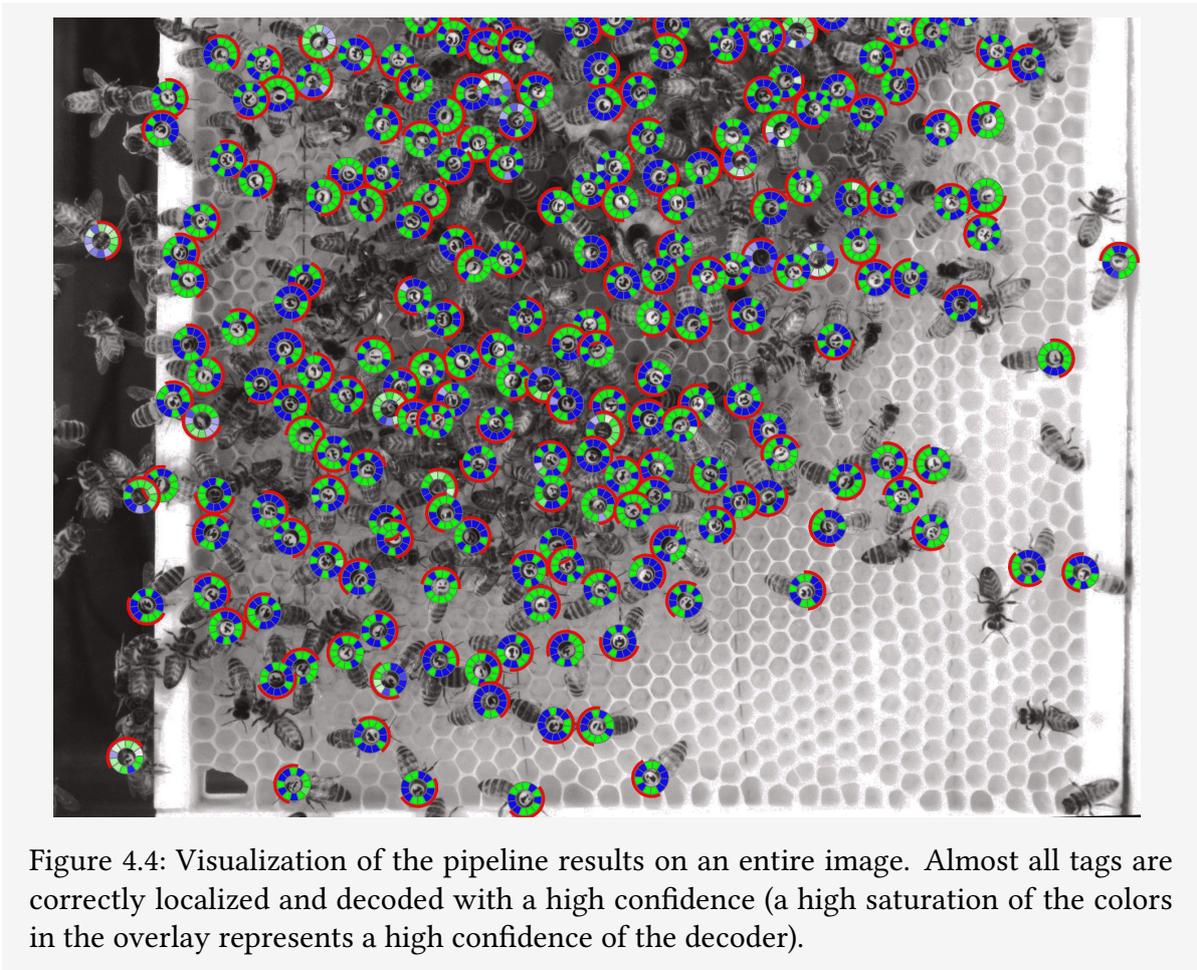


Figure 4.4: Visualization of the pipeline results on an entire image. Almost all tags are correctly localized and decoded with a high confidence (a high saturation of the colors in the overlay represents a high confidence of the decoder).

In total, 3.614.742.669 honeybees were detected on 67.972.617 images for the 2015 recording season (≈ 53 detections per image). In 2016, there were 6.331.078.577 detections in 59.680.181 images (≈ 106 detections per image). The serialized results are 274 GB in size for 2015 and 476 GB for 2016.

Jobs on the HLRN Cray supercomputer are billed in NPL (north-german parallel computer work units). We used Intel Xeon Haswell compute nodes for which 0.1 NPL are billed per core hour. All in all, processing both datasets was billed at 39.896 NPL which corresponds to roughly 398.954 core hours in total. Under optimal circumstances with no IO contention, the data could therefore be processed in roughly one week using 1200 CPUs (100 mpp2 compute nodes on the Cray X30 supercomputer). Please note that the data can also be processed in realtime using a consumer GPU (Geforce GTX 1080 Ti) instead of CPUs. In contrast, processing the data of a single recording season using computer vision pipeline described in Wario et al. (2015) was billed with more than 200.000 NPL.

Discussion

We have presented a new method for the automatic localization and decoding of honeybee markers using deep convolutional neural networks. We are able to identify virtually all individuals in a honeybee colony over many weeks. Tracking individual honeybees with barcode-like markers has only recently been proposed by several groups (Wario et al. 2015, Crall et al. 2015, Blut et al. 2017, Gernat et al. 2018) and will likely allow unprecedented perspectives on the complex interplay between all colony members. While all other systems rely on planar rectangular markers, we decided for an unconventional curved and round marker design. This decision was motivated by our interest in long-term observations and preliminary experiments showing that planar markers endure less mechanical stress and might fall off after only a few days.

While the image processing steps to detect and decode rectangular markers are well researched, we found that conventional, i.e. non-neural computer vision algorithms were less suitable for our custom markers, especially regarding the runtime. Using deep convolutional neural networks as described in this paper, we reached similar, or better detection and decoding accuracy and significantly improved runtime performance compared to comparable approaches using a more traditional marker design (Blut et al. 2017, Gernat et al. 2018).

The performance comparison in section 4 might be arguable due to differences in methodology and data used. While our system does not use any error correction, we discriminate between detection and decoding. The system used in (Blut et al. 2017) might return a detection only if it could be decoded correctly. In this case, the detection accuracy may actually indicate the decoding accuracy, i.e. the proportion of correctly decoded markers. The authors report variable detection performance (between 90.8% and 98.2%) and attribute lower performance to situations in which many bees are in motion. The general level of activity in a bee colony depends on many factors. Under natural conditions there may be a fair amount of motion throughout the day and we therefore expect the decoding accuracy to vary closer to the lower bound. Observations over night may result in decoding accuracies near the upper bound.

In (Gernat et al. 2018) the authors report 87% mean detection accuracy which is lowest among the systems. It remains unclear how many false positives these detections contain. The decoding accuracy for this system (98.58%) is highest among the systems. It remains unknown how many of the potential false detections can be identified in the decoding step. In the best case, this system, hence, has an overall accuracy of 85.7%. The BeesBook system described here yields a similar combined detection and decoding accuracy of $98.3\% \cdot 87.8\% = 86.3\%$.

The system described in (Crall et al. 2015) was validated without manually generated ground truth data and therefore was excluded from the comparison.

In Boenisch et al. (2018) we propose an additional postprocessing step to significantly increase the decoding accuracy. Linking corresponding detections through time (tracking) and averaging bit probabilities of detections within a path increased the combined detection and decoding accuracy to 96%.

Depending on the research question, different properties may be relevant for the end-user of such tracking systems. The BeesBook system offers a very low rate of false detections, a high decoding accuracy and longevity of the markers with relatively inexpensive hardware components. We think this system is suited to tackle a broad range of research questions and we invite researchers to test it.

Our method relies on two machine learning models, the localization and decoding networks. Because of this design decision, we can adapt our software easily to changes in the recording setup. Significant modifications, for a example a new tag design, requires retraining both models, which in most cases may be much easier than redesigning a pipeline of conventional computer vision steps. In our system, the following steps have to be taken: A new training dataset for the localization task has to be created manually. A new localization model can then be trained. The new model can then be used to generate training data for the RenderGAN (Sixt et al. 2016). Finally, data generated by the RenderGAN can then be used to train a new decoder model. Thus, the only manual labor is clicking sample markers images for the localization training dataset. All remaining steps are automated. In many cases, retraining our models may not even be necessary. Exploratory experiments with cameras used in field experiments suggest that the convolutional neuronal networks are surprisingly invariant to changes in illumination, lens distortion and image quality.

A significant result is the 100-fold speedup we achieved in comparison to our previous prototype. Using a consumer-grade graphics card (Geforce GTX 1080 Ti) we obtained realtime performance with a 3 Hz recording framerate and approximately 800 individuals in the colony. In the past, we have recorded full image datasets spanning up to nine weeks of continuous experiment. In total, we currently store three datasets (63 days and 45 TB each) on tape drives granted by our project partner, the North-German Supercomputing Alliance (HLRN). Sample image data is available upon request. A trajectory dataset of all animals over three continuous days is available online (Boenisch et al. 2017).

References

- Balbuena, M. S., Molinas, J. & Farina, W. M. (2012), 'Honeybee recruitment to scented food sources: correlations between in-hive social interactions and foraging decisions', *Behavioral Ecology and Sociobiology* **66**(3), 445–452. 00011.
URL: <http://dx.doi.org/10.1007/s00265-011-1290-3>
- Becher, M. A., Grimm, V., Thorbek, P., Horn, J., Kennedy, P. J. & Osborne, J. L. (2014), 'Beehave: a systems model of honeybee colony dynamics and foraging to explore multifactorial causes of colony failure', *Journal of Applied Ecology* **51**(2), 470–482.
- Beekman, M., Sumpter, D. J. T., Seraphides, N. & Ratnieks, F. L. W. (2004), 'Comparing foraging behaviour of small and large honey-bee colonies by decoding waggle dances made by foragers', *Functional Ecology* **18**(6), 829–835.
- Bengio, Y., Courville, A. & Vincent, P. (2012), 'Representation Learning: A Review and New Perspectives', *arXiv:1206.5538 [cs]*. arXiv: 1206.5538.
- Beshers, S. N. & Fewell, J. H. (2001), 'Models of division of labor in social insects', *Annual review of entomology* **46**(1), 413–440.
- Blut, C., Crespi, A., Mersch, D., Keller, L., Zhao, L., Kollmann, M., Schellscheidt, B., Fülber, C. & Beye, M. (2017), 'Automated computer-based detection of encounter behaviours in groups of honeybees', *Scientific Reports* **7**(1), 17663.
URL: <https://www.nature.com/articles/s41598-017-17863-4>
- Boenisch, F., Rosemann, B., Wild, B., Dormagen, D., Wario, F. & Landgraf, T. (2018), 'Tracking all members of a honey bee colony over their lifetime (in preparation)'.
DOI: 10.7303/syn11737848.1
- Boenisch, F., Rosemann, B., Wild, B., Wario, F., Dormagen, D. & Landgraf, T. (2017), 'Bees-Book Recording Season 2015 Sample'. synapse.org/#!Synapse:syn11737848
DOI: 10.7303/syn11737848.1
- Bonabeau, E., Dorigo, M. & Theraulaz, G. (1999), *Swarm intelligence: from natural to artificial systems*, number 1, Oxford university press.
- Chou, T. H., Ho, C. S. & Kuo, Y. F. (2015), QR code detection using convolutional neural networks, in '2015 International Conference on Advanced Robotics and Intelligent Systems (ARIS)', pp. 1–5.

- Clevert, D.-A., Unterthiner, T. & Hochreiter, S. (2015), 'Fast and Accurate Deep Network Learning by Exponential Linear Units (ELUs)', *Under review of ICLR2016* (1997), 1–13.
- Crall, J. D., Gravish, N., Mountcastle, A. M. & Combes, S. A. (2015), 'BEEtag: a low-cost, image-based tracking system for the study of animal behavior and locomotion', *PloS one* **10**(9), e0136487.
- De Marco, R. J., Gurevitz, J. M. & Menzel, R. (2008), 'Variability in the encoding of spatial information by dancing bees', *Journal of Experimental Biology* **211**(10), 1635–1644.
- Dornhaus, A., Klügl, F., Oechslein, C., Puppe, F. & Chittka, L. (2006), 'Benefits of recruitment in honey bees: effects of ecology and colony size in an individual-based model', *Behavioral Ecology* **17**(3), 336–344.
- Fiala, M. (2005), ARTag, a fiducial marker system using digital techniques, in '2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)', Vol. 2, IEEE, pp. 590–596. 00730.
- Gernat, T., Rao, V., Middendorf, M., Dankowicz, H., Goldenfeld, N. & Robinson, G. (2018), 'Automated monitoring of behavior reveals bursty interaction patterns and rapid spreading dynamics in honey bee social networks (in press)', *Proceedings of the National Academy of Sciences*.
- Girshick, R., Donahue, J., Darrell, T. & Malik, J. (2013), 'Rich feature hierarchies for accurate object detection and semantic segmentation', *arXiv:1311.2524 [cs]*. arXiv: 1311.2524.
- Glorot, X., Bordes, A. & Bengio, Y. (2011), Deep sparse rectifier neural networks, in 'Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics', pp. 315–323.
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A. & Bengio, Y. (2014), 'Generative Adversarial Networks', *arXiv:1406.2661 [cs, stat]*. arXiv: 1406.2661.
- Goyret, J. & Farina, W. M. (2005), 'Non-random nectar unloading interactions between foragers and their receivers in the honeybee hive', *Naturwissenschaften* **92**(9), 440–443. 00027.
URL: <http://dx.doi.org/10.1007/s00114-005-0016-7>
- Grinchuk, O., Lebedev, V. & Lempitsky, V. (2016), Learnable Visual Markers, in 'Advances in Neural Information Processing Systems', pp. 4143–4151.

- Grüter, C. & Farina, W. M. (2009), 'The honeybee waggle dance: can we follow the steps?', *Trends in Ecology & Evolution* **24**(5), 242–247.
- Grüter, C., Acosta, L. E. & Farina, W. M. (2006), 'Propagation of olfactory information within the honeybee hive', *Behavioral Ecology and Sociobiology* **60**(5), 707–715. 00060.
- Grüter, C. & Farina, W. M. (2009), 'Past Experiences Affect Interaction Patterns Among Foragers and Hive-Mates in Honeybees', *Ethology* **115**(8), 790–797. 00017.
URL: <http://onlinelibrary.wiley.com/doi/10.1111/j.1439-0310.2009.01670.x/abstract>
- Grüter, C. & Ratnieks, F. L. W. (2011), 'Honeybee foragers increase the use of waggle dance information when private information becomes unrewarding', *Animal Behaviour* **81**(5), 949–954. 00000.
URL: <http://dx.doi.org/10.1016/j.anbehav.2011.01.014>
- He, K., Zhang, X., Ren, S. & Sun, J. (2015), 'Deep Residual Learning for Image Recognition', *7*(3), 171–180.
- Ioffe, S. & Szegedy, C. (2015), 'Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift', *Arxiv*.
- Khan, Z., Balch, T. & Dellaert, F. (2004), A rao-blackwellized particle filter for eigentracking, in 'Computer Vision and Pattern Recognition, 2004. CVPR 2004. Proceedings of the 2004 IEEE Computer Society Conference on', Vol. 2, IEEE, pp. II–II. 00195.
- Kingma, D. P. & Welling, M. (2013), 'Auto-encoding variational bayes', *arXiv preprint arXiv:1312.6114*.
- Krizhevsky, A., Sutskever, I. & Hinton, G. E. (2012), Imagenet classification with deep convolutional neural networks, in 'Advances in neural information processing systems', pp. 1097–1105.
- Landgraf, T., Rojas, R., Nguyen, H., Kriegel, F. & Stettin, K. (2011), 'Analysis of the Waggle Dance Motion of Honeybees for the Design of a Biomimetic Honeybee Robot', *PLoS ONE* **6**(8), e21354.
- Landgraf, T., Rojas, R. & others (2007), 'Tracking honey bee dances from sparse optical flow fields'.
- Mersch, D. P., Crespi, A. & Keller, L. (2013), 'Tracking individuals shows spatial fidelity is a key regulator of ant social organization', *Science* **340**(6136), 1090–1093.

- Nieh, J. C. (2010), 'A Negative Feedback Signal That Is Triggered by Peril Curbs Honey Bee Recruitment', *Current Biology* **20**(4), 310–315.
URL: <http://www.sciencedirect.com/science/article/pii/S0960982210000758>
- Pizer, S. M., Amburn, E. P., Austin, J. D., Cromartie, R., Geselowitz, A., Greer, T., Romeny, B. t. H., Zimmerman, J. B. & Zuiderveld, K. (1987), 'Adaptive Histogram Equalization And Its Variants', *Computer vision, graphics, and image processing* **39**(3), 355–368.
- Razavian, A. S., Azizpour, H., Sullivan, J. & Carlsson, S. (2014), 'CNN Features off-the-shelf: an Astounding Baseline for Recognition', *arXiv:1403.6382 [cs]*. arXiv: 1403.6382.
- Richter, M. R. & Waddington, K. D. (1993), 'Past foraging experience influences honey bee dance behaviour', *Animal Behaviour* **46**(1), 123–128. 00060.
URL: <http://www.sciencedirect.com/science/article/pii/S000334728371167X>
- Sandstorm Development Group (2013–2017), 'Cap'n proto serialization/rpc system', <https://capnproto.org/>.
- Scheiner, R., Abramson, C. I., Brodschneider, R., Crailsheim, K., Farina, W. M., Fuchs, S., Grünewald, B., Hahshold, S., Karrer, M., Koeniger, G. et al. (2013), 'Standard methods for behavioural studies of apis mellifera', *Journal of Apicultural Research* **52**(4), 1–58.
- Seeley, T. D. (1995), *The wisdom of the hive: the social physiology of honey bee colonies*, Harvard University Press, Cambridge, Mass.
- Shelhamer, E., Long, J. & Darrell, T. (2016), 'Fully Convolutional Networks for Semantic Segmentation', *arXiv:1605.06211 [cs]*. arXiv: 1605.06211.
- Simonyan, K. & Zisserman, A. (2014), 'Very Deep Convolutional Networks for Large-Scale Image Recognition', *arXiv:1409.1556 [cs]*. arXiv: 1409.1556.
- Sixt, L., Wild, B. & Landgraf, T. (2016), 'RenderGAN: Generating Realistic Labeled Data', *arXiv:1611.01331 [cs]*. arXiv: 1611.01331.
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. (2014), 'Dropout: A simple way to prevent neural networks from overfitting', *The Journal of Machine Learning Research* **15**(1), 1929–1958.
- Sun, C., Shrivastava, A., Singh, S. & Gupta, A. (2017), 'Revisiting Unreasonable Effectiveness of Data in Deep Learning Era', *arXiv:1707.02968 [cs]*. arXiv: 1707.02968.

- Szegedy, C., Ioffe, S. & Vanhoucke, V. (2016), 'Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning', *Arxiv* pp. 12–12.
- Veeraraghavan, A., Chellappa, R. & Srinivasan, M. (2008), 'Shape-and-Behavior Encoded Tracking of Bee Dances', *IEEE Transactions on Pattern Analysis and Machine Intelligence* **30**(3), 463–476.
- Visscher, P. K. & Seeley, T. D. (1982), 'Foraging Strategy of Honeybee Colonies in a Temperate Deciduous Forest', *Ecology* **63**(6), 1790–1801.
- von Frisch, K. (1965), *Tanzsprache und Orientierung der Bienen*, Springer.
- Wario, F., Wild, B., Couvillon, M. J., Rojas, R. & Landgraf, T. (2015), 'Automatic methods for long-term tracking and the detection and decoding of communication dances in honeybees', *Behavioral and Evolutionary Ecology* p. 103.
- Wario, F., Wild, B., Rojas, R. & Landgraf, T. (2017), 'Automatic detection and decoding of honey bee waggle dances', *PloS one* **12**(12), e0188626.