

Self-supervised deep visual servoing for high precision peg-in-hole insertion

Rasmus Laurvig Haugaard, Anders Glent Buch, and Thorbjørn Mosekjær Iversen

Abstract—Many industrial assembly tasks involve peg-in-hole like insertions with sub-millimeter tolerances which are challenging, even in highly calibrated robot cells. Visual servoing can be employed to increase the robustness towards uncertainties in the system, however, state of the art methods either rely on accurate 3D models for synthetic renderings or manual involvement in acquisition of training data. We present a novel self-supervised visual servoing method for high precision peg-in-hole insertion, which is fully automated and does not rely on synthetic data. We demonstrate its applicability for insertion of electronic components into a printed circuit board with tight tolerances. We show that peg-in-hole insertion can be drastically sped up by preceding a robust but slow force-based insertion strategy with our proposed visual servoing method, the configuration of which is fully autonomous.

I. INTRODUCTION

Many industrial insertion tasks require high precision for successful completion. When the insertion tolerances are near or lower than the accumulated uncertainties in the system, naive planning-based insertion becomes unreliable. This can be handled either by decreasing the uncertainties in the system or increasing the tolerances.

For some tasks careful calibration is sufficient to reduce the system uncertainties to within the tolerances. However, since the uncertainties in a system depend on the combined effect of many different factors which are hard to accurately model, there is a practical limit to the accuracy of a system. For insertion tasks with very small tolerances, good calibration is, therefore, often not enough.

A common method for increasing the tolerances of an insertion is to introduce either active or passive compliance perpendicular to the insertion direction. However, during the search, before there is peg-hole contact, there is no force-feedback perpendicular to the insertion direction, and prior to peg-hole contact, the force-feedback is thus limited to the binary feedback of whether the peg has hit the surface surrounding the hole, or has been, at least partially, inserted.

Binary force feedback can be used for robust insertion as part of an exhaustive search within the region of uncertainty around the estimated insertion point, e.g. using spiral search. While this has been shown to work successfully for PCB assembly [1], it is a slow technique if the error is significant compared to the tolerances. Since the area that must be explored increases quadratically with the magnitude of the error, $|e|$, and the area of successful insertion positions

This work was supported by Innovation Fund Denmark through the project MADE FAST.

All authors are from SDU Robotics, Maersk Mc-Kinney Møller Institute, University of Southern Denmark.

{r.l.h.a., anbu., t.h.m.i.}@mmmi.sdu.dk

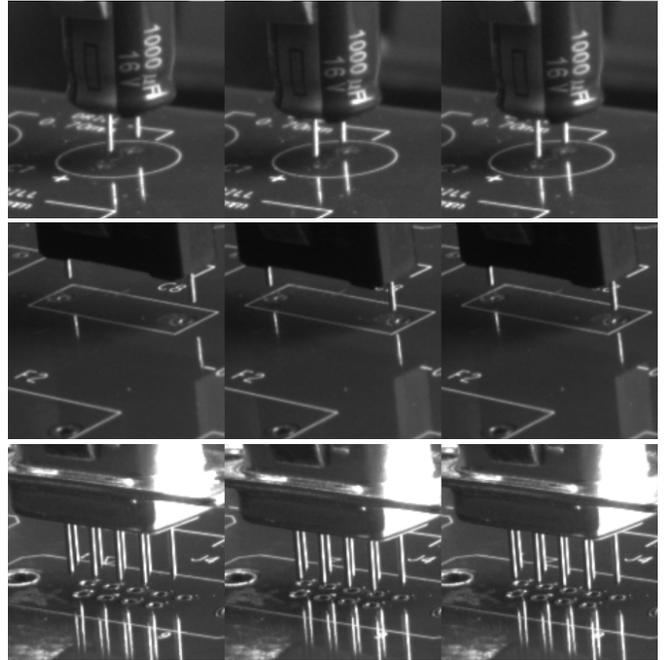


Fig. 1. Examples of two iterations of visual servoing on three PCB components. From left to right: before visual servoing, after first iteration, after second iteration.

increase quadratically with the tolerance, ϵ , the time for insertion, using an exhaustive search, is proportional to the squared ratio between the error and tolerance,

$$t \propto \left[\frac{|e|}{\epsilon} \right]^2. \quad (1)$$

Consequently, the insertion time becomes prohibitively large when tolerances are low relative to the system uncertainties.

Visual servoing is a technique by which a robot is controlled to reduce the error e based on visual feedback. We argue that visual servoing and exhaustive search should be combined to leverage the robustness of slow force-feedback, and visual servoing's ability to reduce the error, to obtain fast and robust insertions.

The literature on visual servoing reports that deep learning based methods can be successfully applied to peg-in-hole insertion. However, the methods rely either on representative 3D models for synthetic renderings or manual involvement in the acquisition of annotated training images. The reliance on accurate models or manual configuration limits the practical use and scalability of such methods.

We present a visual servoing method which is fully

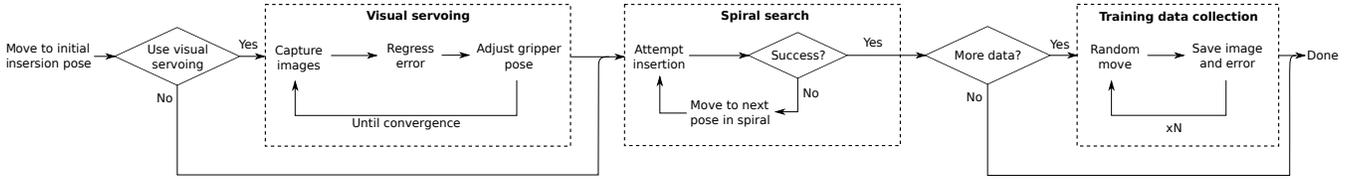


Fig. 2. Overview of proposed system. For the initial insertions, visual servoing is skipped, and spiral search is used to gather training data autonomously. Concurrently, while the system is running, models are trained on the gathered data, and the validation error informs when to enable visual servoing.

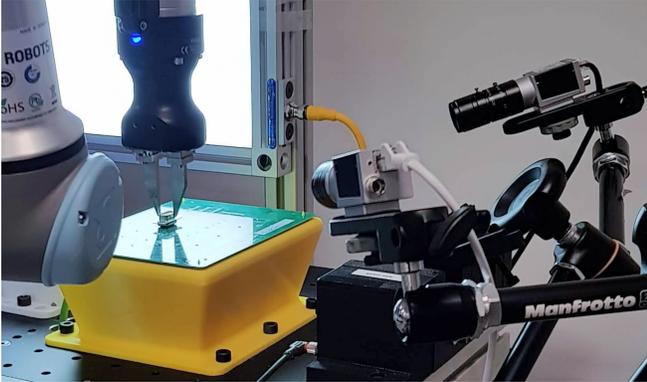


Fig. 3. The experimental setup with a robot, a TCP-mounted gripper, a PCB fixture, and two cameras. The specific light source is designed for another process in the cell and is not necessary for our method.

automated and does not require a model of the object. The proposed method constrains the servoing task to an in-plane 2D alignment of the peg relative to the hole. Our method uses the slow but robust strategy of spiral search to drive self-supervised learning of in-plane alignment by regression of the alignment error in each camera of a multi camera setup. In our experimental setup, we use two cameras. See Fig. 3.

Our contributions are:

- an in-plane deep learning based visual servoing method for high precision peg-in-hole insertion
- and a self-supervised learning scheme based on spiral search which fully automates the configuration and learning process with no reliance on a 3D model.

The proposed method is evaluated for the insertion of five different types of electronic components into a PCB. The evaluation shows that our proposed method enables significantly faster peg-in-hole insertion with sub-millimeter precision.

The paper is organized as follows: Section II discusses the recent literature related to high precision deep visual servoing. Section III then presents the proposed method. This is followed by section IV, which presents the experimental setup, the evaluation procedure, and the results. Finally, V presents the conclusion of the work.

II. RELATED WORK

Visual servoing is defined as a technique by which a robot is controlled by visual feedback [2]. A control loop determines the movement of a robot’s end effector such that an error function, based on visual input, is minimized. The visual input comes from vision sensors which are either mounted on the robot (eye-in-hand) or next to the robot (eye-to-hand). Visual servoing techniques are traditionally categorized depending on the domain in which the error function is defined. Image based visual servoing (IBVS) defines the error function in image space while position based visual servoing (PBVS) defines the error function in euclidian space [3].

Most IBVS techniques fall into one of two categories. The first is feature based methods, which defines the error function in terms of the difference between extracted and desired visual features such as keypoints [4]. The second is direct visual servoing (DVS), in which the error function is not defined explicitly on extracted features but instead defined based on the difference between the current image and a target image [5].

Recent advances in visual servoing rely on deep learning. While a part of the visual servoing literature focuses on achieving large basins of convergence or high generalizability (e.g. [6], [7]), the following focuses on literature in which the main objective is the ability to handle tight tolerances for grasping or insertion.

There are several recent methods which use deep learning for DVS. One such method is [8] which uses a convolutional auto-encoder to extract a learnt latent space representation of current and target image. The robot is then controlled using a derived interaction matrix, which relates difference in current and target latent vector to a desired change in robot motion. Other works uses siamese networks to extract latent features followed by either a network which regresses the velocity vector of the camera [9] or the relative pose between the current and desired image [10], [11]. The evaluation presented in [10] demonstrates that successful insertion of a tool mounted male sub-D connector (peg) into a female sub-D connector (hole) can be achieved using this scheme. However, unlike our method, they fixate the peg to the robot TCP and only perform servoing with respect to the hole, effectively assuming that the only error is in the hole pose. This approach requires low uncertainty on the peg’s in-

hand pose which cannot be guaranteed in general for tasks where objects are grasped by an end effector rather than rigidly mounted on the robot. Furthermore, their acquisition of training data requires initial manual guidance of the robot to the ground truth pose, in order to obtain a target image and initialize the acquisition of training images. Note that the need for a target image is at the core of DVS techniques. This assumption about a single target image is problematic in cases where within-class variations render the information provided by a single target insufficient. In case of electrical components for example, a single target image does not provide information about whether to align with respect to the pins or the body of the component, when the pins are slightly bent. Our method does not assume a single target image and can thus capture such information from multiple targets.

Visual servoing for peg-in-hole insertion has also been done using keypoints extracted with deep learning. Two recent works show that training on synthetic images, combined with heavy domain randomization, can bridge the sim to real gap sufficiently for peg in hole insertions on selected cases. In [12] two cameras are mounted in an eye-in-hand setup. The center of a hole and the tip of a tool mounted cylindrical peg are regressed using convolutional neural networks (CNN), and from these points the euclidian error is computed and used to control the robot end effector. While the peg in hole task is successfully performed, it is only demonstrated to work on insertion of simple cylindrical objects, matching the geometry of the peg model used in the synthetic data. In [13], two networks are used: an encoder-decoder network for self-supervised learning of keypoint extraction and a fully connected network for regressing required robot motion from keypoints extracted from an eye-in-hand stereo camera setup. However, the insertion tasks have high tolerances and inherent mechanical compliance to help guide the insertion. Unlike our method, both [12] and [13] train on synthetic data and thus rely on 3D models and successful sim to real transfer. While the sim to real transfer has been shown to work on selected cases, it depends on the quality of the 3D models and the synthetic renders in general and it is thus not trivial to estimate whether or not it will work for a given task, since the error on even a synthetic validation dataset does not necessarily represent the performance on real data. Also, when sim to real transfer fails, expert knowledge is required to improve the system. In contrast, our training and validation data is drawn from the target distribution, and we are thus able to leverage standard machine learning practices for reliable model evaluation before deployment.

To the best of our knowledge, our work is the first visual servoing method to demonstrate sub-millimeter precision peg-in-hole insertion of non-trivial objects, requiring neither 3D models for synthetic image generation nor manual involvement in the acquisition of training images.

III. METHOD

This section first defines and formalizes the in-plane visual servoing task, then presents our method, including how we

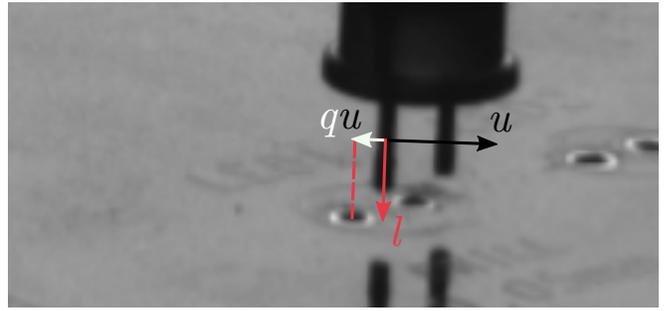


Fig. 4. The insertion direction, l , is shown with a red arrow, the error direction, u , is shown in black, and qu is shown in white, all projected into the image for visualization.

formulate in-plane visual servoing as a well-posed learning problem, the visual servoing control loop, and the strategy for self-supervision.

A. In-plane visual servoing

We define in-plane visual servoing as the control of a robot where the movement of the end effector is constrained to a constant rotation and a displacement confined to a plane, that is $\frac{d\phi}{dt} = 0$ and $\frac{dp}{dt} \cdot l = 0$, where $\phi \in \text{SO}(3)$ is the rotation, $p \in \mathbb{R}^3$ is the TCP position, and $l \in \mathbb{R}^3$ is the insertion direction and thus a normal vector to the alignment plane.

Constraining the visual servoing task to positional, 2-DOF, alignment is a simplification compared to full 5-DOF alignment, followed by insertion along the last DOF. However, small errors in the 3-DOF rotation can be seen as a reduction of the insertion tolerances, and orientation constraints from mechanical feeders and grippers help to reduce these errors. We thus argue that the simplification is valid in many cases, and hypothesize that the simplification leads to faster and more robust visual servoing.

B. Proposed visual servoing method

The prerequisite for our method is a calibrated setup, like the one seen in Fig. 3, with a robot, a gripper attached to the robot TCP and two or more cameras. The cameras in our experimental setup are attached to the table, but can be attached in-hand or to a separate robot for increased flexibility. Our method relies on system calibration to extract crops and to relate the insertion direction to the images, but the visual feedback loop ensures that the method is robust towards calibration uncertainties.

The reason for choosing at least two cameras is to create a well posed regression problem. The in-plane alignment error could theoretically be regressed from a single image by relative depth estimation from perspective effects. However, since the objects are small compared to the distance between camera and insertion, the camera projection is close to orthographic in the region of interest. It could also be argued that the in-plane alignment of the peg can be learned from the peg's projection onto the image plane relative to the projection of the background. However, this would introduce strict requirements for the distance between hole plane and peg to be constant, which in turn would reduce the robustness

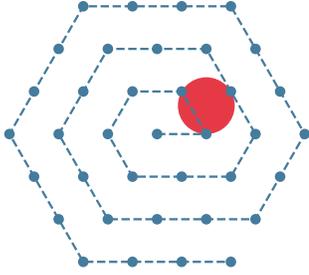


Fig. 5. Our in-plane spiral search pattern in blue. The largest enclosed circle is drawn in red. The circle diameter represents the minimum tolerance for which the pattern is guaranteed to succeed.

with respect to uncertainties, including but not limited to grasping uncertainties, part variations and time-dependent calibration errors. Consequently, we argue that refraining from regressing the depth alignment and instead obtain the in-plane alignment from triangulation from two or more cameras, as in [12], is a significantly better posed problem. We hypothesize that a well-posed servoing problem is key to learning a mapping that is robust to uncertainties and generalizes well to new instances.

Given an insertion direction, $l \in \mathbb{R}^3$, the position of a camera, $c \in \mathbb{R}^3$, and the approximate position of the insertion, $x \in \mathbb{R}^3$, we define the view vector, $v = x - c$, and the camera specific error direction, $u = \frac{l \times v}{|l \times v|}$. Given an in-plane alignment error, $e \in \mathbb{R}^3$, $e \cdot l = 0$, the scalar error along the error direction is $q = e \cdot u$. See Fig 4. Under the assumption of orthographic projection in the region of interest, the error can be reconstructed, like in [12], as the least squares solution to the set of linear equations established by scalar errors from multiple views,

$$\begin{bmatrix} u_{1x} & u_{1y} & u_{1z} \\ u_{2x} & u_{2y} & u_{2z} \\ \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} e_x \\ e_y \\ e_z \end{bmatrix} = \begin{bmatrix} q_1 \\ q_2 \\ \vdots \end{bmatrix}. \quad (2)$$

We propose to learn a mapping, $f_\theta : \mathbb{R}^{r \times r} \mapsto \mathbb{R}$, from a camera image, I , with resolution r to the scalar error, q , however, to make the dynamic range of the output independent of object scale and image resolution, we learn the mapping to $y = qf(rz)^{-1}$ instead, where f is the focal length in pixels, z is the approximate depth of the insertion in the camera, and $y \in \mathbb{R}$ is the normalized scalar error in the image. Specifically, we aim to learn the parameters, θ of a CNN, f_θ , that minimizes the mean squared error loss

$$L = \frac{1}{N} \sum_i^N (y_i - f_\theta(I_i))^2. \quad (3)$$

During inference, q can be established from y . The full visual servoing control loop is shown in Algorithm 1.

C. Autonomous configuration and continuous learning

The visual servoing network presented in Sec. III-B could in principle be trained using manually annotated images

Algorithm 1: In-plane visual servoing

Input: Insertion direction, l . Approximate hole position, h . Approximate camera positions, c_j and focal lengths, f_j , of n_{cams} cameras. Number of iterations, n_{iters} .

$p \leftarrow$ current TCP position

for $i = 1, \dots, n_{\text{iters}}$ **do**

for $j = 1, \dots, n_{\text{cams}}$ **do**

$I_j \leftarrow$ capture image from camera j

$y_j \leftarrow f_\theta(I_j)$ // est. norm. error

$v_j \leftarrow h - c_j$ // view vector

$u_j \leftarrow \frac{l \times v_j}{|l \times v_j|}$ // error direction

$q_j \leftarrow y_j r z f_j^{-1}$ // error magnitude

end

$(U, q) \leftarrow \left(\begin{bmatrix} u_{1x} & u_{1y} & u_{1z} \\ u_{2x} & u_{2y} & u_{2z} \\ \vdots & \vdots & \vdots \end{bmatrix}, \begin{bmatrix} q_1 \\ q_2 \\ \vdots \end{bmatrix} \right)$

$\hat{e} \leftarrow$ solve $U\hat{e} = q$ by least squares

$p \leftarrow p + \hat{e}$

 move TCP to p

end

and/or synthetic images. However, the strength of our method lies in having autonomous collection of training data from the target distribution. This allows the configuration of the visual servoing method for previously unseen objects to be reliable and fully autonomous. See Fig. 2.

Peg-in-hole insertion requires poses of peg and hole, and depending on the system, this is approximately known based on calibration, demonstration, pose estimation, vibrational feeders, and/or by other means. Assuming that the tolerance for insertion is lower than the uncertainty on the relative hole to peg pose, a search strategy is needed to ensure successful insertion, and when such a method is in place, like spiral search, the insertion is guaranteed given enough time.

When an insertion has been successfully completed, we know that the end TCP position led to a successful insertion, and thus we can assume that to be the correct in-plane position. Since most industrial robots have a low *relative* positional uncertainty, it is possible to acquire annotated training data by capturing images with *known* in-plane errors. From such a dataset it is straight forward to compute the ground truth scalar error, y , in each image.

We can relax the previous assumption that the in-plane position at a correct insertion position is *the* correct, fully centered position. The mean squared error loss assumes that $p(y|I)$ is normally distributed, and leads to regressing the mean of said distribution. Even though the insertions in the training data are not fully centered, the vision model should still learn to regress the center, as long as the assumption that $p(y|I)$ is normally distributed around the true y is good.

The autonomously obtained dataset is split into a training and a validation set, and the model is trained with early

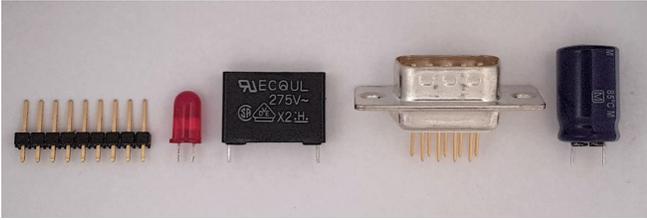


Fig. 6. The five electronic components used in the experiments. We refer to them, from left to right, as PH, LED, C1, DSUB and C2.

stopping, monitoring the loss on the validation set to inform when to stop to avoid overfitting. We split the data such that all data points from an insertion either goes to the training or the validation set to enable early stopping that specifically monitors the model’s ability to generalize to new insertions, in contrast to sampling the validation set randomly from all data points. The performance on the validation set can be used to decide whether the model should be deployed on the system, or if more data is required to obtain a good model.

During deployment, the performance of the model can also be monitored to detect a distributional shift which could be caused by changes in the objects, e.g. due to supplier-change, or by changes to other parts of the system, like a new background, lighting, time-dependent calibration errors, etc. The system would then be able to collect more data and autonomously obtain robustness to the variations that occur in the system.

IV. EXPERIMENTS

We evaluate our visual servoing method on insertion of electronic components into a PCB. The experimental setup is shown in Fig. 3, and we choose five different types of components, shown in Fig. 6, covering variance in size, number of pins and visual appearance. The insertion tasks have sub-millimeter tolerances, lower than the accumulated errors in system calibration, grasp uncertainties, and PCB tolerances, making it a relevant case for our method. The setup has mechanical fixtures and linear vibrational feeders, from which the robot can grasp the components.

As introduced in Sec. I, spiral search provides a robust but slow way to deal with the fact that the uncertainties are larger than the insertion tolerances. A common spiral search implementation lets the peg slide on the surface surrounding the hole while applying a force in the insertion direction, registering if the peg dips into the surface, indicating an at least partial insertion. In our case however, a such approach could bend the pins and damage the PCB surface. Instead, our spiral search implementation attempts insertions in a spiral-like pattern at the intersections of an isometric grid, as shown in Fig. 5. This, like the common archimedean spiral search, provides guarantees for the diameter of the largest enclosed circle, the minimum allowed tolerance, and is more efficient than a regular grid.

As discussed in Sec. III-C, we use the robust but slow search to obtain visual servoing training data autonomously. See Fig. 2. Specifically, we perform ten insertions for each

TABLE I
MEAN INSERTION TIMES IN SECONDS
WITH AND WITHOUT VISUAL SERVOING.

	PH	LED	C1	DSUB	C2	Avg
vs	1.7	1.8	2.0	1.7	2.3	1.9
no vs	11.8	47.5	42.7	11.3	47.8	32.2

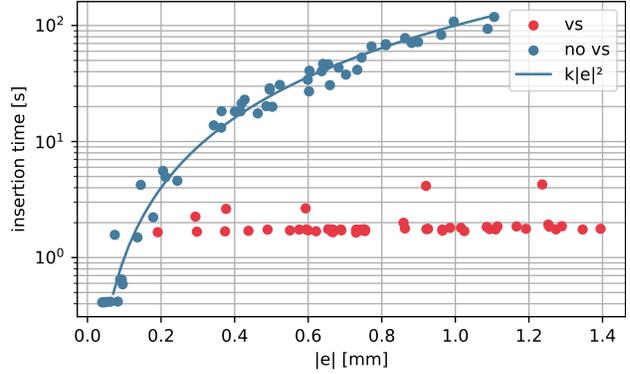


Fig. 7. Insertion times as a function of retrospective error. Note the logarithmic scale on the time axis. k is a constant.

component type and capture images at 100 sampled TCP positions per insertion. The TCP positions are sampled around the successful in-plane position. We sample the direction of the error along the plane uniformly, the magnitude along the direction uniformly from 0 to 1 mm, and an offset perpendicular to the plane, also uniformly between 0 and 1 mm.

For each component type, a vision model is trained to predict y , as discussed in Sec. III-B. Any capable vision architecture can be used. In our experiments, we use a ResNet50 [14] pretrained on ImageNet [15] classification, replacing the last fully connected layer with a two-layer MLP for regressing y . The MLP has 128 neurons and ReLU activations at the pre-ultimate layer. The models are trained with the Adam [16] optimizer and a learning rate of 10^{-4} . We use the data from eight of the insertions for training and two for validation to enable early stopping that specifically monitors the model’s ability to generalize to new insertions.

The total configuration time per component was approximately 25 minutes, including 10 minutes for data collection and 15 minutes for training on an RTX2080 GPU.

To evaluate our method’s robustness and speed with respect to uncertainties, we sample errors uniformly from an in-plane disc with a radius of 1 mm, in addition to the inherent system uncertainties. We then attempt spiral search insertion with and without preceding visual servoing. For both approaches we attempt ten insertions per component, amounting to 100 insertions in total. We use three iterations of visual servoing and a spiral search tolerance of 0.1 mm for all insertions.

The insertion times are presented in Table I, where insertion time refers to the duration from when the TCP is immediately above the hole in the sampled error position

with the grasped component, to when a successful insertion is determined by spiral search. We perform the spiral search with up to 1 mm error, which is the size of the error we sample during the experiments, however, since the error is added to the accumulated errors inherent in the system, 2 of the 50 spiral search insertions without visual servoing fails. All the insertions with visual servoing succeeds, while being significantly faster across all component types and more than 15 times faster on average compared to insertion without visual servoing.

Because the sampled errors are added to the system uncertainties, we do not know the actual error a priori, however, like we obtain the dataset, we can use the position at successful insertion to estimate the initial error, retrospectively. The relationship between the retrospective initial error and the insertion time is visualized in Fig. 7 for all 100 insertions. Pure spiral search is fast, when the error is very small, but increases quadratically with the error. In contrast, the insertion time with visual servoing is approximately constant.

Note that the retrospective initial errors are estimates and are only exact up to the actual insertion tolerances which are different between the five component types, and largest for PH and DSUB, which also shows in Table I. Insertion based solely on spiral search will tend to find the closest successful insertion, explaining the generally lower retrospective errors without visual servoing.

The critical error, where visual servoing outperforms pure spiral search, depends on the tolerances of a given task. Fig. 7 indicates that visual servoing leads to a reduction in insertion time when the accumulated system uncertainty is more than approximately 0.2 mm.

41 of 50 insertions with visual servoing inserts directly at the center point of the spiral search, compared to 5 of 50 insertions without visual servoing. The average retrospective error after visual servoing is 0.03 mm.

Of the average insertion time of 1.9 s with visual servoing, 1.3 s is spent on visual servoing itself and the remaining 0.6 s is spent on the spiral search. Of the 1.3 s spent on visual servoing, 0.5 s is spent on capturing images, 0.4 s is spent on forward passes on the vision model, and 0.4 s is spent on physically moving the robot. We use rather long exposure times, capture the images sequentially, and run the model inference on a laptop CPU. Adding a bright light source, capturing the images in parallel and running inference on a GPU could thus reduce the time spent on image acquisition and inference.

Also, continuous visual servoing, as proposed in [12] would be able to further reduce the impact of image acquisition time, because images are captured in parallel with robot motion, and reduce the time for robot motion, since the robot does not need to come to a stop for image acquisition. Note however, that the increase in speed comes with added implementation complexity in terms of image timestamp synchronization, and the visual servoing path's dependence on acquisition and inference speed.

V. CONCLUSION

This paper presented a novel self-supervised deep visual servoing method for high precision peg-in-hole insertion. The method is fully automated and does not rely on the availability of 3D models. This is achieved by constraining the visual servoing task to in-plane alignment and training a convolutional neural network to regress scalar alignment errors in image space in a dual camera setup. Annotated data is gathered autonomously using a robust but slow force-based search method. The method has been evaluated on insertion of PCB components. The evaluation showed that preceding a robust but slow search strategy with our proposed method reduced the average insertion time by an order of magnitude.

REFERENCES

- [1] S. F. Mathiesen, L. C. Sørensen, D. Kraft, F. Hagelskjær, and T. M. Iversen, "Towards flexible pcb assembly using simulation-based optimization," in *Towards Sustainable Customization: Bridging Smart Products and Manufacturing Systems: CARV 2021, MCPC 2021*. Springer, 2021, pp. 166–173.
- [2] B. Siciliano, L. Sciacivco, L. Villani, and G. Oriolo, *Force control*. Springer, 2009.
- [3] B. Thuilot, P. Martinet, L. Cordesses, and J. Gallice, "Position based visual servoing: keeping the object in the field of vision," in *Proceedings 2002 IEEE International Conference on Robotics and Automation (Cat. No. 02CH37292)*, vol. 2. IEEE, 2002, pp. 1624–1629.
- [4] L. Weiss, A. Sanderson, and C. Neuman, "Dynamic sensor-based control of robots with visual feedback," *IEEE Journal on Robotics and Automation*, vol. 3, no. 5, pp. 404–417, 1987.
- [5] C. Collewet, E. Marchand, and F. Chaumette, "Visual servoing set free from image processing," in *2008 IEEE International Conference on Robotics and Automation*. IEEE, 2008, pp. 81–86.
- [6] Y. Harish, H. Pandya, A. Gaud, S. Terupally, S. Shankar, and K. M. Krishna, "Dfvs: Deep flow guided scene agnostic image based visual servoing," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 9000–9006.
- [7] F. Sadeghi, A. Toshev, E. Jang, and S. Levine, "Sim2real viewpoint invariant visual servoing by recurrent control," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4691–4699.
- [8] S. Felton, P. Brault, E. Fromont, and E. Marchand, "Visual servoing in autoencoder latent space," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 3234–3241, 2022.
- [9] S. Felton, E. Fromont, and E. Marchand, "Siame-se (3): regression in se (3) for end-to-end visual servoing," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 14454–14460.
- [10] C. Yu, Z. Cai, H. Pham, and Q.-C. Pham, "Siamese convolutional neural network for sub-millimeter-accurate camera pose estimation and visual servoing," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2019, pp. 935–941.
- [11] F. Tokuda, S. Arai, and K. Kosuge, "Convolutional neural network-based visual servoing for eye-to-hand manipulator," *IEEE Access*, vol. 9, pp. 91820–91835, 2021.
- [12] R. Haugaard, J. Langaa, C. Sloth, and A. Buch, "Fast robust peg-in-hole insertion with continuous visual servoing," in *Conference on Robot Learning*. PMLR, 2021, pp. 1696–1705.
- [13] E. Y. Puang, K. P. Tee, and W. Jing, "Kovis: Keypoint-based visual servoing with zero-shot sim-to-real transfer for robotics manipulation," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 7527–7533.
- [14] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [15] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [16] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *ICLR (Poster)*, 2015.