

Deep Paper Gestalt

Jia-Bin Huang
Virginia Tech
jbh Huang@vt.edu

Abstract

Recent years have witnessed a significant increase in the number of paper submissions to computer vision conferences. The sheer volume of paper submissions and the insufficient number of competent reviewers cause a considerable burden for the current peer review system. In this paper, we learn a classifier to predict whether a paper should be accepted or rejected based solely on the visual appearance of the paper (i.e., the gestalt of a paper). Experimental results show that our classifier can safely reject 50% of the bad papers while wrongly reject only 0.4% of the good papers, and thus dramatically reduce the workload of the reviewers. We also provide tools for providing suggestions to authors so that they can improve the gestalt of their papers.

1. Introduction

Peer review — a thorough examination of a scholarly work by other experts in the community — is an essential aspect of disseminating scientific results. However, the record-breaking number of paper submissions to top-tier computer vision conferences and the insufficient number of competent reviewers make the peer review process increasingly more difficult (see Figure 1). To review all these submissions, conference organizers have to expand the pool of reviewers and inevitably include less experienced students [3]. Consequently, the authors who spent months or years of efforts on a paper submission may end up receiving poorly justified, ill-considered, or unfair reviews.

In this paper, we address this pressing issue in two aspects. First, we train a deep convolutional neural network using prior conference proceedings to determine the quality of the paper based on its visual appearance (known as paper gestalt [19]). Second, we provide diagnostic tools to help authors enhance their future paper submissions. Trained on ICCV/CVPR conference and workshop papers from 2013 - 2017, our deep network based classifier achieves 92% accuracy on papers in CVPR 2018. Our model safely rejects the number of bad paper submissions by 50% while sacri-

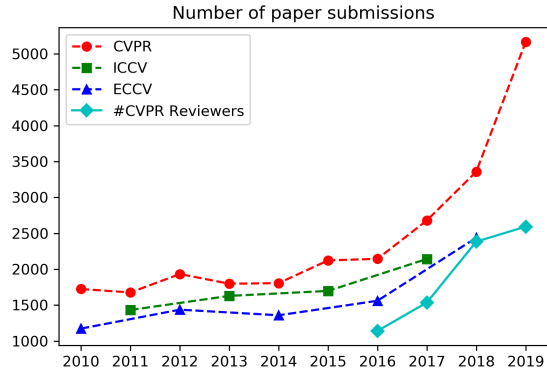


Figure 1. **Need.** The number of the paper submissions to top-tier computer vision conferences has been increased dramatically over the past few years. The number of competent reviewers (as shown in the cyan curve), however, does not grow nearly as fast.

ficing only 0.4% of good paper submissions. Our system can thus be used as a *pre-filter* in a cascade of the paper review process. Using the collected Computer Vision Paper Gestalt (CVPG) dataset, we can 1) visualize class-specific discriminative regions of a particular paper submission or 2) translate a bad paper to a good one directly. These tools help inform the authors *where* and *how* to improve the gestalt of their papers.

2. Related Work

Administrative methods. Several methods have been proposed to address the surge in the number of paper submissions through administrative policies. Examples include desk-reject by area/program chairs (e.g., violation of anonymity, formatting, or clearly out of scope), mandatory abstract submission one week before the paper submission deadline, expansion of the reviewer pool, and training materials for inexperienced reviewers [1].

Text-based methods. Automatic grading techniques have been developed for grading essay [14], response to mathe-

matical questions [13], and handwritten work [18]. These techniques, however, do not take into account the rich visual information available in the paper and may be subject to bias toward popular keywords trending in the community.

Our tool for improving paper gestalt is related to sentence editing [9, 21] and automatic random paper generator for computer science [4] and math [2]. Our approach differs in that we directly learn the mapping in the image space.

Vision-based methods. Computer vision techniques have been applied to accessing the quality of actions [17], surgical skills [26], and images [20, 23]. The work most related to our work is that of the awesome Bearnensquash [19], where the AdaBoost algorithm is used for learning the good/bad paper classifier. Building upon the methodology in [19] that relies on hand-crafted visual features, we revisit the paper gestalt problem with deep learning and learn task-specific representation through an end-to-end training process.

3. Learning to Recognize Good/Bad Papers

We leverage deep convolutional neural networks (ConvNets) to learn discriminative representation based solely on the visual appearance of a paper, known as *paper gestalt*. In the following, we start with describing the problem formulation and presenting our dataset construction process. We then provide the implementation details of the network training. We validate the performance through an empirical evaluation on a held-out testing set and visualize the class-specific discriminative regions produced by the trained network.

3.1. Problem formulation

We formulate the problem as a binary classification task. Our training dataset consists of N labeled data samples, $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$, where x_i denotes the i -th paper and $y_i \in \{0, 1\}$ is binary label indicating whether the i^{th} paper x_i is a good paper or a bad one. Our goal here is to learn a function $F_\theta(\cdot)$ parametrized by θ that can recognize good/bad papers from unseen paper submissions (e.g., paper submissions to future conferences).

3.2. Dataset construction

Data source. We collect positive examples (good papers) from the list of accepted papers in top-tier computer vision conferences. Specifically, we gather the Open Access versions of the accepted papers from recent conferences sponsored by the [Computer Vision Foundation \(CVF\)](#). This includes six CVPR and three ICCV proceedings from 2013 to 2018.

For negative examples, as we do not have access to papers that were rejected from these conferences, we follow

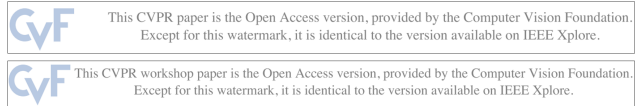


Figure 2. **Preventing data leakage.** The header of the Open Access versions of the papers contains the information that the classification models try to predict (i.e., good or bad papers). While the differences are small in the resized images, our results show that our classification network can easily achieve 100% accuracy on both the training and testing set when the headers were not removed, suggesting that the network found a way to “cheat”.

[19] and use workshop papers as an approximation. Similar to the conference papers, we gather the Open Access versions of all the workshop papers from the CVF website. Note that these negative examples can be noisy as some of the papers 1) were also accepted at the main conferences or 2) were not submitted to the main conference. At the same time, these workshop papers can also be viewed as “hard negative” examples as many of the papers have been significantly improved by addressing the comments from the reviewers.

Data acquisition and preprocessing. Here we outline the detailed steps we used for constructing the dataset.

1. **Crawl:** We crawl both the positive and negative examples from the [CVF Open Access website](#).
2. **Filter:** As some of the workshop papers have different page limits as the main conference papers (e.g., 6 pages including references), the classification task becomes trivial for these cases. We therefore keep only papers with sufficient (≥ 7) pages.
3. **PDF2Image:** We use the [pdf2image](#), a python wrapper for [pdftoppm](#), to convert the downloaded PDFs to images. We arrange these pages into a 2×4 grid. For papers with a missing 8^{th} page, we pad it with a blank page. We also discard the pages greater than 8 (mainly references cited in the paper). The original size of the converted image is of size 2200×3400 pixels.
4. **Pre-processing:** To prevent the data leakage problem, we remove the header on top of the first page. Without this preprocessing step, the learned classifier can become overly optimistic or even invalid because the classifier can focus on the header region while ignoring the visual contents of the paper.

The detailed statistics of the collected Computer Vision Paper Gestalt (CVPG) dataset are shown in Table 1. There are in total 5618 positive examples and 1503 negative examples. Figure 3 shows random samples from both positive



Figure 3. **Random samples of the collected Computer Vision Paper Gestalt (CVPG) datasets.** Glancing through samples in both classes show that there are differences in terms of the general layout of the paper. Our goal here is to leverage deep ConvNets to learn representation for capturing these patterns.

Table 1. **Computer Vision Paper Gestalt (CVPG) dataset.**

Positive examples		Negative examples	
Venue	# samples	Venue	# samples
CVPR 2013	471	CVPR-W 2013	80
ICCV 2013	454	ICCV-W 2013	101
CVPR 2014	540	CVPR-W 2014	61
CVPR 2015	602	CVPR-W 2015	113
ICCV 2015	526	ICCV-W 2015	116
CVPR 2016	643	CVPR-W 2016	184
ICCV 2017	621	ICCV-W 2017	350
CVPR 2017	783	CVPR-W 2017	251
CVPR 2018	978	CVPR-W 2018	247
Total	5618	Total	1503

and negative samples of the collected dataset. The dataset is available on our project website <https://github.com/vt-vl-lab/paper-gestalt>.

3.3. Paper review as image classification

To simulate the actual potential usage of our system (i.e., predicting good/bad papers from unseen paper submissions), we use the positive/negative examples in the CVPR 2018 as our testing set and the papers in the prior conferences/workshops from 2013 to 2017 as our training set.

We use ResNet-18 [10] (pre-trained on ImageNet) as our classification network.¹ We replace the ImageNet 1,000 class classification head with two output nodes (good or bad papers). Following the practice of transfer learning, we finetune the ImageNet pre-trained network on the proposed CVPG dataset with stochastic gradient descent (SGD) with a momentum of 0.9 for a total of 50 epochs. We set the initial learning rate as 0.001 and decay it by

¹Deeper networks with larger capacity can also be used. However, we do not observe significant performance improvement when using ResNet-34 or ResNet-50.

a factor of 0.1 every ten epochs. To accommodate the class-imbalanced training data, we use the weighted cross-entropy loss (weighted by the inverse of the training examples in each class). We resize all the images to 224×224 pixels for both training and testing. We choose not to apply standard data augmentation techniques such as random cropping, horizontal flipping, or photometric transformation during training to keep the original visual content and layout of the entire paper. The network training process takes less than 30 minutes on a NVIDIA Titan V100 GPU.

3.4. Experimental results

Evaluation. On the test dataset (CVPR 2018 conference/workshop papers), our trained network achieves an overall accuracy of 92%. By varying the threshold on the network predictions after the softmax layer, we plot the ROC curve to further characterize the performance of our model in Figure 4. The x-axis shows the *false positive rate (FPR)*, indicating the portion of bad papers getting accepted by our model. The y-axis shows the *false negative rate (FNR)*, indicating the portion of good papers getting rejected by our model. Note that here we plot the FNR instead of the true positive rate (TPR) to better illustrate the trade-off between the two error types.

Choosing different threshold values leads to different trade-off. For example, if we allow only 0.4% of the good papers getting rejected, we can accurately reject 50% of the bad ones. If we allow 5% of the good papers getting rejected (as there will be inevitable noises in peer reviews anyway), we can reject up to 75% of the bad papers.

Here we use a more concrete example to better understand what the results mean. There are in total 3309 valid submissions to CVPR 2018 with 979 of them accepted (good papers) and 2230 papers rejected (bad papers). Assuming the actual negative examples show the same distributions in the workshop papers, applying our model (with 0.4% FPR and 50% FNR) to all the valid submissions to

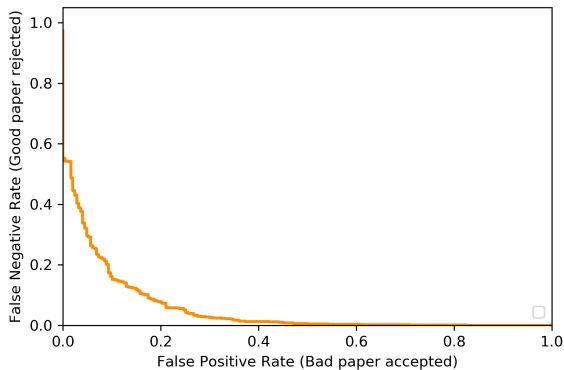


Figure 4. **Performance characterization of the trained good/bad paper classifier.** The x-axis denotes the false positive rate (the percentage of bad papers getting accepted). The y-axis denotes the false negative rate (the percentage of good papers getting rejected).

CVPR 2018 can safely reject 1115 bad papers (without peer reviews) at the cost of sacrificing 4 good papers (among the 979 good ones). Such an automatic pre-filtering stage substantially reduces the workload of reviewers.

Class-specific activation maps. While the model achieves decent classification performance, we believe that it is unlikely that the classifier will ever be used in an actual conference. Nevertheless, we can take a closer look at how the classification model makes the decision and in turns improve the paper gestalt of our future paper submissions.

Multiple visualization techniques for understanding deep neural networks have been proposed. For examples, retrieving image patches that maximize a particular neuron [7], reconstructing the input image [16], quantifying the interpretability of latent layers [5], and mapping the activations to the input image space with a deconvNet [22]. In this paper, we use the class-specific activation mapping [24] for visualizing the discriminative regions for the classifying a paper into a good or a bad one.

Figure 5 shows sample class-specific activation maps on papers that were rejected by our model. In the first row, the discriminative regions generated by our classifier highlight mostly on the *incomplete pages*. This makes sense because well-polished papers often squeeze the contents compactly into precisely 8 pages (with some clever use of `vspace`). In the second row, the class-specific activation maps focus on the top-right corner of the first page. It appears that the classifier picks up the *absence of a motivation or teaser figure* in the first page as its primary reason for rejecting a paper. This reveals that it is crucial to include a motivation figure on the first page to illustrate the main idea of the

work.

Figure 6, on the other hand, shows the class-specific activation maps on papers that were accepted by our model. The discriminative regions for good papers include the teaser figure (easier to understand the main idea), detailed tables (comprehensive experimental evaluation), and colorful images (qualitative results). We believe that such visualization can be applied as a diagnostic tool to help identify the strength/weakness of one’s paper submissions in the future.

Self-evaluation. Following [19], we also convert this paper as an image as shown in Figure 7. We then apply our trained classification network to determine whether this paper should be accepted or rejected. Unfortunately, despite the visually pleasing figures/tables/plots in the paper, our classifier predicts a posterior probability of 97.4% that this paper is a bad one and should be rejected. We attribute the primary weakness of our paper to the incomplete pages.

4. Learning to Enhance Paper Gestalt

In addition to classifying a paper and highlighting discriminative regions, here we aim to provide further suggestions to help authors enhance the paper gestalt of their submissions.

4.1. What does a good paper look like?

One approach for providing suggestions is to generate visual layouts of a good paper. To this end, we train a *good paper generator* using generative adversarial networks (GANs) [8]. Specifically, we train our generator using the state-of-the-art progressively growing GANs [12]. We use the conference papers from 2013-2017 as our training dataset. The entire training process takes about a week with two NVIDIA Titan V100 GPUs.

Figure 8 shows 15 random samples generated by our trained model. We see that these synthesized good papers often have a balanced layout of figures/tables/plots/equations. However, the visual quality of these samples is poor, particularly on the generated blurry “figures” and “tables”. The results are expected because every figure/table in the training dataset is unique.

For more examples, please see the latent space interpolation video on <https://youtu.be/yQLsZLf02yg>.

4.2. Learning bad \rightarrow good paper translation

While the generated layout by the *good paper generator* looks convincing, its practical usage is limited as it is difficult to follow a particular template when preparing a paper submission. Instead, we often wish to *translate* a bad paper into a good one. As we do not have access to pairs

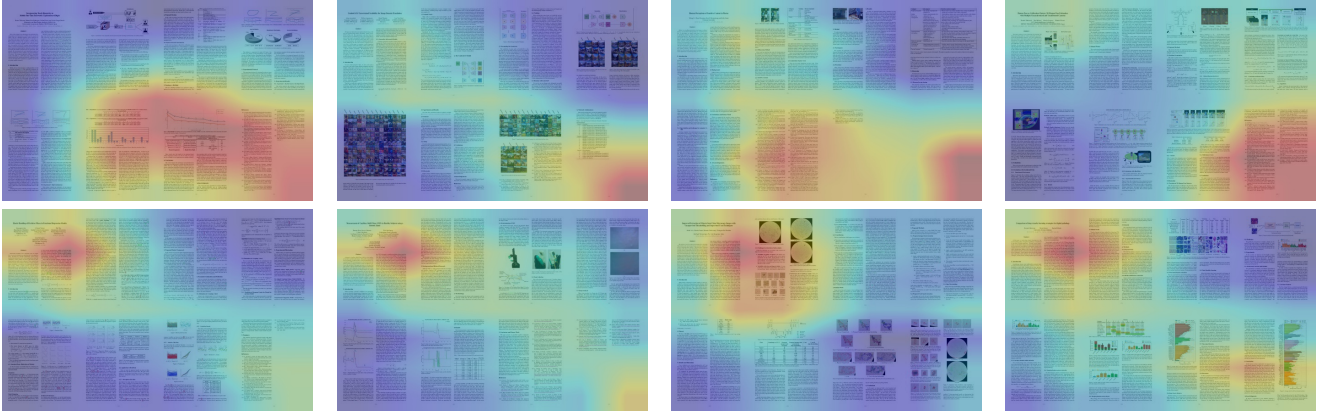


Figure 5. **Class-specific discriminative regions for bad papers.** (Top) Failing to fill the paper into a full eight-page paper is a discriminative visual cue for bad paper. (Bottom) The generated heatmaps focus on the top-right corner of the first page. This suggests that the *absence* of illustrative figures in the first two pages may cause the paper more difficult to understand.



Figure 6. **Class-specific discriminative regions for good papers.** The heatmap generated by class activation mapping [24] highlights regions specific to good papers, e.g., teaser figures in the first page for illustrating the main ideas, tables/plots showing a sense of thoroughness in experimental validation, impressive math equations, and arrays of colorful images for qualitative results from benchmark datasets.

of good/bad papers, we use the cycle-consistent adversarial network for unpaired image-to-image translation [25]. Similar to the experimental setting in training a classifier, we use the conference/workshop papers in 2013-2017 for training and the papers in CVPR 2018 for testing. Figure 9 shows the animation flipping between the input image (a randomly selected workshop paper) and its translated version (best viewed using the Adobe Reader). Our model automatically suggests several changes for the input papers. For examples, adding the teaser figure (for clarity), adding figures at the last page (usually failure cases), filling the incomplete last page, and making the figures more colorful.

5. Limitations and Discussions

In this paper, we revisit the problem of paper gestalt using modern deep ConvNets. We show that the model trained on existing paper proceedings (e.g., ICCV/CVPR

2013-2017) generalizes well to unseen paper submissions (e.g., CVPR 2018). Our classifier safely rejects the number of bad paper submissions by half while only sacrificing 0.4% of the good paper submission. Our classifier can thus serve as an effective pre-filter to significantly reduce reviewers’ workload in the later stages. Our model also runs very fast, takes a only a few seconds to classify thousands of paper submissions. In addition to automatic determining a paper submission should be accepted or rejected, we also introduce several diagnostic tools (visualizing discriminative regions, a good paper generator, and a bad-to-good paper translation network) that can help authors improve the gestalt of their papers.

While interesting results have been shown, our work suffers from the following limitations. First, our classifier relies entirely on the visual appearance of a paper. Ignoring the actual paper contents may wrongly reject papers with



Figure 7. **This paper.** We apply the trained classifier to this paper. Our network ruthlessly predicts with high probability (over 97%) that this paper should be rejected without peer review. ☹

good materials but bad visual layout or accept crappy papers with good layout. Second, both our classifier and the generative model assume that all the papers have the same typesetting style (provided by the conference template). As a result, the trained classifier cannot be applied to other conference papers with different formatting styles. This limits the applicability of our method because other related fields (ML/NLP/AI/Robotics) also experience similar issues of rapidly increased workload in paper review. One potential solution to this is to convert papers from image space to a high-level abstraction, e.g., a structured representation of text blocks, figures, and tables (e.g., using the method in [6]). Such an abstraction is therefore invariant to various typesetting styles required by different venues. Third, the bad-to-good paper generator (trained with [25]) can only produce one single output as a good paper. To generate diverse paper editing suggestions, we may use recent methods such as [15, 11]. Fourth, the collected training samples can be very noisy (mainly for the negative examples) because we do not have access to the rejected papers at the main conference. The new OpenReview model adopted by the International Conference on Learning Representation (ICLR) offers a way to gather *ground truth* positive and negative samples for training paper gestalt classifier.

References

[1] How to write good reviews for cvpr. <https://www.dropbox.com/s/725p60wcajbb8xh/How%20to%20Review%20for%20CVPR.pptx?dl=0>. 1

[2] Mathgen. <http://thatismathematics.com/mathgen/>. 2

[3] Nips 2018: How do i write a good review? https://www.reddit.com/r/MachineLearning/comments/8ite3n/r_nips_2018_how_do_i_write_a_good_review/. 1

[4] SCIGen - an automatic cs paper generator. <https://pdos.csail.mit.edu/archive/scigen/>. 2

[5] D. Bau, B. Zhou, A. Khosla, A. Oliva, and A. Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *CVPR*, 2017. 4

[6] C. A. Clark and S. K. Divvala. Looking beyond text: Extracting figures, tables and captions from computer science papers. In *AAAI Workshop: Scholarly Big Data*, 2015. 6

[7] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014. 4

[8] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014. 4

[9] K. Guu, T. B. Hashimoto, Y. Oren, and P. Liang. Generating sentences by editing prototypes. *Transactions of the Association of Computational Linguistics*, 6:437–450, 2018. 2

[10] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 3

[11] X. Huang, M.-Y. Liu, S. Belongie, and J. Kautz. Multimodal unsupervised image-to-image translation. In *ECCV*, 2018. 6

[12] T. Karras, T. Aila, S. Laine, and J. Lehtinen. Progressive growing of gans for improved quality, stability, and variation. In *ICLR*, 2018. 4

[13] A. S. Lan, D. Vats, A. E. Waters, and R. G. Baraniuk. Mathematical language processing: Automatic grading and feedback for open response mathematical questions. In *Proceedings of the Second ACM Conference on Learning@ Scale*, 2015. 2

[14] L. S. Larkey. Automatic essay grading using text categorization techniques. In *International ACM SIGIR conference on Research and development in information retrieval*, 1998. 1

[15] H.-Y. Lee, H.-Y. Tseng, J.-B. Huang, M. Singh, and M.-H. Yang. Diverse image-to-image translation via disentangled representations. In *ECCV*, 2018. 6

[16] A. Mahendran and A. Vedaldi. Understanding deep image representations by inverting them. In *CVPR*, 2015. 4

[17] H. Pirsiavash, C. Vondrick, and A. Torralba. Assessing the quality of actions. In *ECCV*, 2014. 2

[18] A. Singh, S. Karayev, K. Gutowski, and P. Abbeel. Gradescope: a fast, flexible, and fair system for scalable assessment of handwritten work. In *Proceedings of the Fourth ACM Conference on Learning@ Scale*, 2017. 2

[19] C. von Bearnensquash. Paper gestalt. In *Secret Proceedings of Computer Vision and Pattern Recognition*, 2010. 1, 2, 4

[20] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on image processing*, 13(4):600–612, 2004. 2

[21] J. Weston, E. Dinan, and A. H. Miller. Retrieve and refine: Improved sequence generation models for dialogue. *arXiv preprint arXiv:1808.04776*, 2018. 2

[22] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *ECCV*, 2014. 4

[23] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 2



Figure 8. **Randomly generated samples of good papers.** These random samples capture the *gestalt* of a good paper: illustrative figures upfront, colorful images, a balanced layout of texts/math/tables/plots.

Figure 9. **Paper enhancement using CycleGAN [25].** The trained bad-to-good paper model can be used as a suggestive tool for translating a bad paper into a good one. Typical suggestions include adding teaser figure upfront, making the figures more colorful, and filling up the last page so that it appears like a well-polished paper. This figure contains *animated images* flipping back and forth between the original bad paper and the translated good paper (best viewed using Adobe Acrobat Reader).

- [24] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba. Learning deep features for discriminative localization. In *CVPR*, 2016. 4, 5
- [25] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017. 5, 6, 7
- [26] A. Zia, Y. Sharma, V. Bettadapura, E. L. Sarin, T. Ploetz, M. A. Clements, and I. Essa. Automated video-based assessment of surgical skills for training and evaluation in medical schools. *International journal of computer assisted radiology and surgery*, 11(9):1623–1636, 2016. 2