

Assisting human experts in the interpretation of their visual process: A case study on assessing copper surface adhesive potency

Tristan Hascoet
Kobe University

tristan@people.kobe-u.ac.jp

Xuejiao Deng
Kobe University

dengxuejiao1005@yahoo.co.jp

Kiyoto Tai
MEC Co., Ltd.

tai295@mec-np.com

Yuji Adachi
MEC Co., Ltd.

Sachiko Nakamura
MEC Co., Ltd.

Tomoko Hayashi
MEC Co., Ltd.

Mari Sugiyama
MEC Co., Ltd.

Yasuo Ariki
Kobe University

Tetusya Takiguchi
Kobe University

Abstract

Deep Neural Networks are often thought to lack interpretability due to the distributed nature of their internal representations. In contrast, humans can generally justify, in natural language, for their answer to a visual question with simple common sense reasoning. However, human introspection abilities have their own limits as one often struggles to justify for the recognition process behind our lowest level feature recognition ability: for instance, it is difficult to precisely explain why a given texture seems more characteristic of the surface of a finger nail rather than a plastic bottle. In this paper, we showcase an application in which deep learning models can actually help human experts justify for their own low-level visual recognition process: We study the problem of assessing the adhesive potency of copper sheets from microscopic pictures of their surface. Although highly trained material experts are able to qualitatively assess the surface adhesive potency, they are often unable to precisely justify for their decision process. We present a model that, under careful design considerations, is able to provide visual clues for human experts to understand and justify for their own recognition process. Not only can our model assist human experts in their interpretation of the surface characteristics, we show how this model can be used to test different hypothesis of the copper surface response to different manufacturing processes.

1. Introduction

Humans are experts in communicating the reasoning process behind their answer to visual questions. For instance, on typical Visual Question Answering (VQA) samples, human annotators are often able to convincingly jus-

tify, in natural language, the reason behind their answer to a certain visual question using simple common sense reasoning. In contrast, deep Learning models are often viewed as black box predictors lacking interpretability in the sense that existing tools often fail to explain the decision process behind the models predictions. For instance, a deep learning model trained end-to-end on a VQA dataset may be able to provide the same answer as its human counterpart, but the process through which the model reaches this answer is entirely opaque.

While it is true that humans can justify for their answers on high level reasoning tasks, humans also often fail to explain the process behind their low-level feature recognition ability: For example, precisely defining the nature of a specific texture (what are the defining features of a plastic or a wooden surface?) or specific low-level part attributes exhibiting large intra-class variations (what is the defining features of a “leg” or a “wing”?) is a very difficult task. Humans constantly perform such low-level visual recognition tasks while being unable to precisely justify for their own recognition process.

In this paper, we present one very practical instance of such a situation in the Printed Circuit Boards (PCB) industry, in which expert material scientists are tasked with assessing the adhesive potency of copper surfaces. We propose a model that, under careful design considerations, is able to provide visual clues for human experts to understand and justify for their decision process.

The proposed model is designed so that a subset of its internal representations carry semantically meaningful information that can be visualized and easily interpreted by humans. Providing these visual clues, however, comes with the cost of imposing additional constraints on the architecture, which we found to degrade the model accuracy: In-

deed, we found that networks with unrestricted architectures, (which do not provide interpretable features) perform better than network architectures restricted so as to provide semantically meaningful representations. This is because, as we restrict the architecture of the model, we formulate an assumption on the impact of the manufacturing process on the surface statistics which may not hold in reality. This result suggests an inherent trade-off between expressivity and the explainability in designing model architecture.

While the degradation of the model accuracy is problematic from a performance perspective, it offers an interesting opportunity from an explainability perspective: As the model accuracy degrades due to the inadequacy of the assumption made by the model, we can use the model accuracy as a proxy metric for the adequacy of different assumptions. This allows us to quantitatively assess different assumptions regarding the impact of manufacturing processes on the copper surface. This may prove useful to quantify the impact of manufacturing process on copper surface adhesive potency and eventually help optimize the manufacturing process.

In essence, the argument this paper is aiming for is as follows: although deep learning models lack the “common sense reasoning” abilities of humans, and the powerful formalism of natural language to communicate and justify for their decision process, they can provide useful tools to visualize and explain low-level recognition processes.

In practice, the contribution of this paper is as follows:

1. We formalize a segmentation procedure based on a probabilistic weak label segmentation framework.
2. We introduce a formalism to show how the model accuracy can be used as a proxy metric to quantify the validity of different assumptions on the dataset.

The remainder of this paper is organized as follows: In Section 2, we present some background information on the motivation for this project: We start by discussing the importance of copper surface adhesive potency, and detail the dataset used in our experiments. Section 3 details our contribution. Section 4 briefly relates our work to different research topics and Section 5 presents the results of our experiments. Finally, Section 6 further discusses the relevance of our results, insisting on the limitations of our assumptions to conclude this paper.

2. Background

Printed circuit boards (PCBs) are an integral part of a wide variety of electronic devices, including industrial and household appliances (e.g. TV and PC), mobile communication devices and automobiles. PCBs play an important role in electrical connection between electronic components. Copper has been used in the PCBs industry as the

conducting material, and the electric copper circuits are isolated from each other by insulators (solder resist, prepreg etc). A multilayered PCBs have a laminate having a plurality of electro-conductive layers with insulating layers interposed therebetween. So there are many interfaces related to the copper and resins in PCBs. Figure 1 illustrates the organization of such an electric circuit.

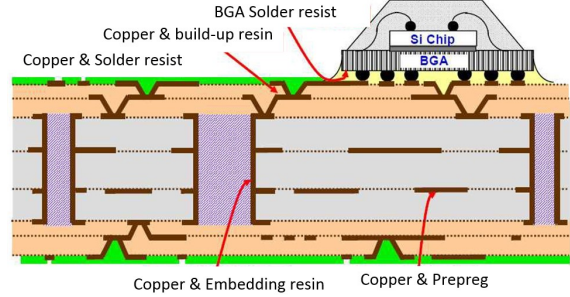


Figure 1. Illustration of a Printed Circuit Board. Copper circuits are made with various insulators for several purposes.

Since PCBs have been required to have higher heat-resistant properties in recent years, copper surface treatment technologies have performed a more and more important role in the manufacturing process. This is because they can enable PCBs to maintain high copper to resin adhesion even under harsh conditions. Copper surface treatment (copper surface roughening) offers one of the most effective ways to increase adhesion between copper and resin. Copper surface roughening has been widely used for the purpose of increasing adhesion of copper to resins. It produces a unique surface topography which enhances the mechanical bonding of copper to resins, as illustrated in Figure 2.



Figure 2. Illustration of a copper roughness surface. (Left) a perfectly smooth surface provides small adhesive surface as the interface between copper, in brown, and the resin, in green is minimal. (Right) a rough surface provides a larger surface at the interface of the resin. Larger contact surface areas provide higher adhesive potency.

In very broad terms, rough surface allow for stronger bonds as the asperities of the surface provide a wide range of adhesive surface area and an anchoring effect. In contrast, smooth surfaces dont provide such effects so that they have lower adhesive potential. In the remaining of this paper, we will refer to the potential bonding strength of a copper surface as its “adhesive potency”. It is also important to note that the adhesive potency of a copper surface is related to its “roughness”, which is observable at the microscopic scale.

Electronic substrate manufacturers have developed advanced manufacturing processes to shape the surface of copper sheets in order to increase their adhesive potency. This is typically achieved by applying an etching solution on the copper surface. Being able to accurately assess the adhesive potency of a copper surface would allow to further optimize manufacturing processes to increase the reliability of electronic devices. However, assessing the adhesive potency of a copper surface is a complex task, even for the most expert practitioners. Hence the motivations of this study is two fold: First we aim to automate the evaluation of a copper sheet adhesive potency from microscopic imaging of its surface. Second, we aim to better understand the affect of different manufacturing processes on the adhesive potency of copper surfaces. Towards this goal, we built a dataset of microscopic images of copper surfaces, which we detail below.

We imaged copper surfaces using Scanning Electron Microscopy (SEM) at a resolution of 100 nm. To investigate the impact of different manufacturing processes on the copper surface, we applied 16 different etching solutions, with decreasing etching power, to the copper surface. For each of these solutions, we captured 50 SEM images of 960×1280 pixels so that the full dataset consists of 800 (16×50) images. Each image is annotated with a label t corresponding to the etching solution used to shape the copper surface. Each solution was obtained by submitting the original solution $t = 0$ to an extreme stress test for a period of time t . Hence, we know that for all images of copper surfaces with label t show higher adhesive potency than the images labelled with $t' > t$. However, we do not know the *exact* impact of the stress test on the surface adhesive potency.

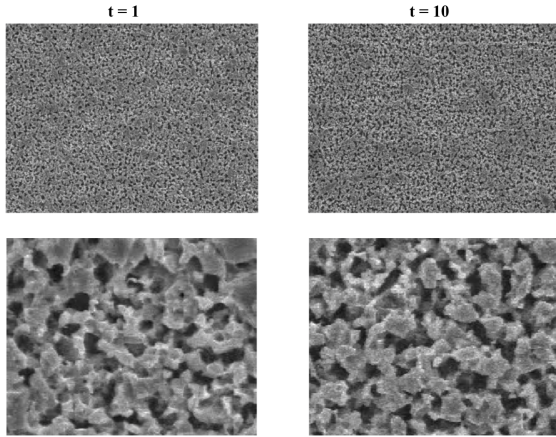


Figure 3. Illustration of a few images from the dataset. (Top) Full images. (Down) Zoomed-in areas of 200×200 pixels. (Right) Sample image of label $t = 1$. (Left) Sample image of label $t = 10$. The difference between both images are minimal to an untrained eye. Precisely defining the visible difference with words is a difficult task.

3. Method

3.1. Dataset and Notations

We denote the dataset described above as $\mathcal{D} = \{(x_i, t_i) | i \in [1, 800]\}$ where x denote gray scale images $x \in \mathbb{R}^{H \times W}$ and labels $t \in [0, 16]$ correspond to the time of stress test applied. We split the dataset \mathcal{D} into a training \mathcal{D}_{tr} , validation \mathcal{D}_{val} and test \mathcal{D}_{te} set so that the number of images x per label t in each set are 40 for training set, and 5 for the validation and test sets.

3.2. Baseline

We start by establishing a strong baseline for our study. The baseline architecture follows standard convolutional network designs for image classification. This architecture, illustrated in Figure 4, is made of several residual blocks sequentially interleaved with max pooling operations. Each residual block consists of n repetitions of a sequence of 3×3 Convolution, Batch Normalization and ReLU layers, followed by a residual skip connection. We set N residual blocks between every max pooling layer and we denote by d the number of pooling layers. Hence, the full depth D of the network (in number of convolution layers) is given by $D = d \times n \times N + 1$, where the term 1 corresponds to the initial 3×3 Convolution layer happening before the first pooling operation. Finally, the top layer of the network is made of a global average pooling layer followed by a linear softmax layer with output dimension 15 corresponding to our number of classes. For simplicity and contrary to standard practices, we keep the number of channels c constant in all layers of the network.

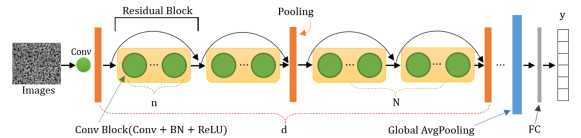


Figure 4. Illustration of our baseline architecture. With our modular architecture definition, the architecture is fully defined by parameters n, N, c and d .

With this parameterization, our network is fully specified by the four hyper parameters c, d, n and N . We performed a grid search over these hyper parameters to select the best performing architecture. The details of this architecture search are given in the experiment section, and, as we shall see then, the best performing architecture performs significantly better than human experts. However, the decision process through which this model reaches such a high accuracy is entirely opaque as the distributed nature of the model's internal representations provides little interpretability. The remainder of this section details our attempt to design an architecture that can provide useful explanations of the process through which high accuracy recognition can be

performed.

3.3. Assumptions

The basic idea behind our method is similar to the visual insights provided by the visualization method of [?] and the attention mechanisms in visual models for image captioning and question answering tasks: We would like to evaluate the contribution of each input pixel to the final classification decision. Indeed, explicitly pointing out the surface regions responsible for low and high adhesive potency would provide guidance for human experts to identify visual patterns characteristic of either cases.

To provide such attention map, we modify our initial problem formulation into a segmentation task: Given an input $x \in \mathbb{R}^{H \times W}$, we want to design a model that outputs a segmentation mask $s \in \mathbb{R}^{H \times W}$ assessing the contribution of each individual pixel to the output adhesive potency score (i.e., the output class t). Training a typical segmentation model for this task would require ground truth segmentation masks s for each image x of the training set. However, manually annotating ground truth segmentation masks for this task is not feasible, as human experts are not able to provide such fine-grained annotations. Instead, we have to train the segmentation given a single ground truth label t per image. In order to do so, we make the following assumption:

Assumption: For a given image x with label t , the *pixel-wise* values of the true (unknown) segmentation mask s follow a spatially stationary binomial distribution whose expected value, averaged over the spatial dimensions of x , is given by t following:

$$s_{hw} \sim \mathcal{B}(f(t)), \forall h, w, t \in H \times W \times T \quad (1)$$

in which we introduced a target function f , which we shall discuss in Section 3.6.

In other words, we suppose that there exists a true binary segmentation map s assigning to each individual pixel of x a binary adhesive potency score: s takes 0 values in regions of the surface providing low adhesive potency (i.e. smooth copper surface areas) and 1 values in regions of the surface providing high adhesive potency (i.e. rough copper surface areas). The adhesive potency score of an entire image x is thus given by the average of the pixel-wise values, and this average value is uniquely defined by the image label t .

Figure 5 provides a visual illustration of this idea. For each image, we suppose the existence of such binary mask s quantifying the adhesive potency of local areas of the surface. The average value of the binary mask s amounts to the ratio of the copper surface covered in white (i.e. with value 1, corresponding to high adhesive potency). Our assumption means that this ratio stays constant for different images x sharing a similar label t , and the exact value of this ratio is given by the target function f .

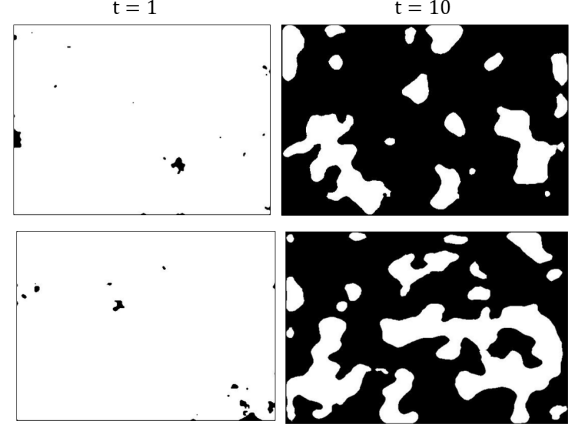


Figure 5. Illustration of binary segmentation masks. White pixels represent areas of high adhesive potency and black surface represent areas of low adhesive potency. We make the assumption that the ratio of white surface is constant for different samples (top and bottom) with equal label t . The exact value of this ratio is defined by the target function f .

3.4. Architecture

In this section, we present the architecture used to compute binary segmentation masks from input images. Our architecture extends the baseline architecture presented in Section 3.1, with an ascending path that progressively up-samples the output of the descending path, similar to the UNet architecture [3] and illustrated in Figure 6.

Residual modules of the ascending path are the exact symmetric of their analog in descending path. The ascending path uses bilinear upsampling layers instead of the max pooling layers of the descending path. Different from the UNet architecture, and following previous works [?], we merge the outputs of the descending path modules with the inputs of the ascending path by summation, instead of concatenation. We also use valid convolutions to preserve the spatial resolution of the output. Finally, we add a sigmoid layer at the top of the network in order to bound the output values between 0 and 1.

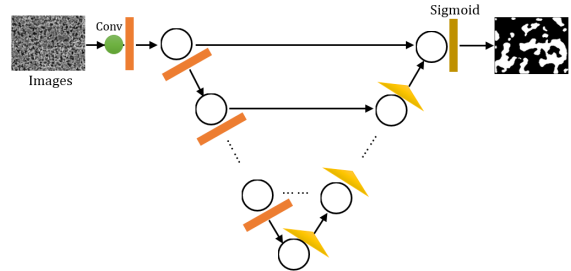


Figure 6. Illustration of the segmentation model architecture. This architecture follows standard practice in UNet-like segmentation architectures.

3.5. Loss Function

Given a segmentation model M_θ , with weight parameters θ , and an input image x , we denote the average of the model output by $m_\theta(x)$:

$$M_\theta : \mathbb{R}^{H \times W} \rightarrow [0, 1]^{H \times W} \quad (2a)$$

$$\tilde{s} = M_\theta(x) \quad (2b)$$

$$m_\theta : \mathbb{R}^{H \times W} \rightarrow [0, 1] \quad (2c)$$

$$m_\theta(x) = \frac{1}{HW} \sum_{h,w} \tilde{s}_{hw} \quad (2d)$$

We can then train the model by regressing the average value of the segmentation mask to the target label given by $f(t)$. Given a training dataset of labeled samples $\mathcal{D}_{tr} = (x_i, t_i)$, learning is thus done by minimizing the following loss function over the model's parameters θ :

$$\theta^* = \operatorname{argmin}_\theta \mathbb{E}_{x,t \in \mathcal{D}_{tr}} \|m_\theta(x) - f(t)\|^2 \quad (3)$$

However, as we shall see in the next section, the target function f represents an unknown ideal function, so we do not have access to the actual values of $f(t)$. Instead, we will approximate f with a known hypothesis function $g \approx f$, so that the actual training loss used in our experiments is:

$$\theta^* = \operatorname{argmin}_\theta \mathbb{E}_{x,t \in \mathcal{D}_{tr}} \|m_\theta(x) - g(t)\|^2 \quad (4)$$

3.6. Target Function

In section 3.3, we have made the assumption that labels t uniquely define the expected value of the ground-truth binary segmentation masks s through a target function f . In this section, we discuss the role of this target function.

The target function $f(t)$ describes the evolution of the copper surface adhesive potency with time t . More precisely f defines the evolution of the *ratio of adhesive surface area* with time (see Figure 5).

However, f is an ideal, unknown function, of which we have only supposed the existence. We do not know the value taken by $f(t)$ for a given t because we do not know the exact impact of the manufacturing process on the copper's surface characteristics. We thus introduce a known hypothesis function g to approximate the ideal target function f . g expresses our belief of what the values taken by the true function $f(t)$ are. Although we do not know the exact values taken by f , we know several of its characteristics, which we can use to reduce the search space of hypothesis functions g :

In Section 2, we have established that, for all t , copper surfaces with label $t' > t$ should have lower adhesive potency than copper surfaces of label t . Hence, f should be a monotonically decreasing function of time, and we can

thus restrict our search of hypothesis function g to monotonically decreasing function. For example, the simplest of such function would be the linear function, taking linearly decreasing values from 1 to 0, which we analyze in the experiment section.

Second, as f is assumed to describe the true evolution of the copper surface statistics with t , a universal function approximator trained on sufficient samples generated by f should generalize well to unseen samples similarly generated by f . In other words, training a segmentation model with a hypothesis $g \approx f$ should lead to low test errors. On the other hand, we expect a function approximator trained on an hypothesis g orthogonal to f to overfit some spurious correlations of the training set, yielding high generalization error. This idea motivates us to use the error of the model on the held out validation dataset as a proxy measure on the validity of the hypothesis function g .

In the experiment section, we evaluate the model errors with different hand-crafted hypothesis functions g , and show that the model trained on a poorly chosen hypothesis function g (i.e. non monotonically decreasing hypothesis) tend to yield large test errors, confirming our idea that generalization error might be used as proxy metric to quantify the validity of a hypothesis function g . In future work, we aim to jointly learn the hypothesis function g with the model parameters θ to automatically formulate hypothesis

4. Related Work

We identify three different lines of research that share similarities with our study: Explainability of learned visual representations, weak supervision of segmentation models and machine learning applications for material science. Our approach can be seen as a weakly supervised approach to provide material experts with strong supports for interpretable decision, which falls at the intersection of these three research lines. We briefly present some of these works in the following subsections

4.1. CNN Interpretability

The operations performed by deep models are hard to interpret: because of the distributed nature of the model's internal representation, it is difficult to assign a useful meaning to each individual neuron or feature map. This is problematic as model failure cases are very hard to investigate and justify. Hence, an interesting line of work is researching for tools to interpret neural network processes. We can see two approaches in this line of work: one is focused on visualization and the other one focused on generating natural language explanations.

On the side of natural language explainability, Park *et al.* [2] jointly trained to answer and justify for their answers on a visual question answering task. They combine visual

attention maps and natural language generation to bring interpretability to the model’s output.

More related to our work is the line of work focused on visualizing neural network hidden activations Yosinski *et al.* [4] provide two useful tools for visualizing and interpreting neural nets. One is to visualize the activations produced on each layer of a trained convnet as it processes an image or video, and the other enables visualizing features at each layer of a DNN via regularized optimization in image space.

More recently, Carter *et al.* [?] propose an explorable activation atlas of vision model’s learned feature by using feature inversion to visualize millions of activations from an image classification network Their technique provides insights regarding the network’s conceptual representations.

Visual attention [?] also helps in explaining the process through which model’s outputs are computed by providing visual clues as to what regions of the input space contribute the most to the final decision. These have been used to investigate model’s operations on image captioning [?] and VQA [?] tasks, and is similar to the idea behind our segmentation model.

4.2. Weakly Supervised Segmentation

An other point we should pay attention to is different labeling strategies in conventional supervised learning. Typically speaking, it’s often assumed that each instance is associated with one single label. However, there should usually be more than one labels for one instance in real-world tasks, if it is multi-label learning.

Eugene et al. propose a semi-supervised framework that employs image-to-image translation between weak labels. Yuxing Tang et al. build a similarity-based knowledge transfer mode trying to investigate whether knowledge about visual and semantic similarities of object categories can help improve the performance of detectors trained in a weakly supervised setting.

The most related work with ours is Rihuan Ke’s[1] in which they present a semi-supervised learning strategy for segmentation with lazy labels and develop a multi-task learning framework to integrate the instance detection, separation and segmentation within a deep neural network.

4.3. Machine Learning for Material Science

Deep learning has showed remarkable success on many important learning problems in chemistry, drug discovery, biology and materials science. In the field of materials science, deep neural networks have also been receiving increasing attention and have achieved great improvements, for example, in material property prediction and new materials discovery for batteries. CNNs have also been used for defect detection on microscopic images of various material surfaces [?].

Beyond the material sciences, we note a growing interest in applying deep learning techniques for scientific discovery. This perhaps best exemplified by the impressive successes of AlphaFold [?] in protein folding estimation, or the Celeste [?] project which cataloged celestial objects of visible universe.

Our work, while much more modest in its scale, shares the characteristic of using vision models to unravel the underlying principle of material reactions to chemical treatments.

5. Experiments

5.1. Classification results

We start by evaluating the baseline classification model described in Section 3.3 and illustrated in Figure 4. The hyper parameters n , N , d , and c of the classifier architecture were obtained by a grid search within the limits of a 12GB Nvidia GPU memory size . The model was trained on the training dataset for 200 epochs using the Adam optimizer with default parameterization.

We compare the classification outputs of the baseline model to those of an expert material scientist on a blind test. The human expert was given the sample images of the training set to practice, and we evaluated the expert’s answers on the samples of the test set.

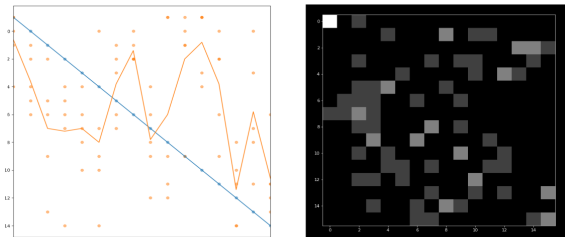


Figure 7. Results of the human evaluation. (Left) Regression plot on the test set. The blue line illustrates the ground-truth labels. Each yellow dot represents the expert’s predicted labels on individual images. The yellow line represents the expert’s average answers for images of similar ground-truth label. (Right) Expert’s answers visualized as a confusion matrix.

Figure 7 shows the results achieved by the model and Figure 8 shows the results achieved by our baseline model. As can be seen in these figures, the model tends to predict the manufacturing process more accurately than the human evaluator. However, these results should be taken with a grain of caution as the human expert was given little time to practice on this specific dataset, while the model was selected as the best performing baseline from an extensive parameter search. We plan on re-conducting the human evaluation in an updated version of this paper.

Interestingly, however, the human expert seems to accurately recognize the surfaces of highest adhesive potency as

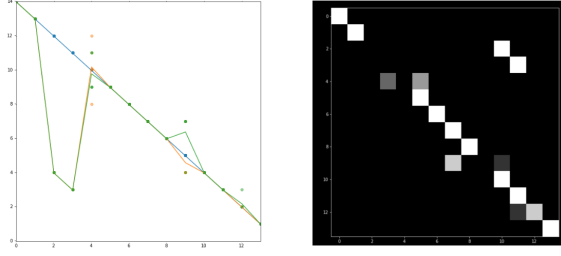


Figure 8. Results of the model evaluation. (Left) Regression plot. Blue dots represent predictions on the training samples, yellow dots represent validation samples and green dots represent the test samples. (Right) Confusion matrix of the test samples

seen in the downward trend of his results for low t values. For example, he easily identified the surface for $t = 0$. On the other hand, his guesses for low adhesive potency are much more random.

This is in stark contrast with the classifier accuracy, which perfectly identifies the low adhesive potency surfaces (high t values), but consistently misclassifies surfaces of high adhesive potency (i.e.; for $t = 2$ and $t = 3$).

5.2. Visualization of Segmentation Results

To motivate our results, Figure 9 illustrates an output of the segmentation model on a small patch of an input image x . To a novice observer, the model seems to assign lower adhesive potential to smoother regions of the input. In future work, we plan on further exploring these visualizations with human experts to better understand the patterns characterizing the surface adhesive potency.

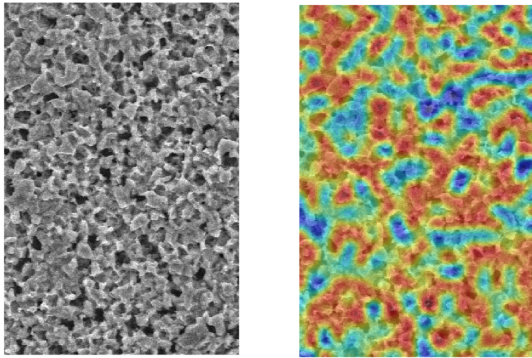


Figure 9. Visualization of the segmentation model output. (Left): Zoomed-in input image. (Right): Model segmentation output overlaid on the input image. Blue areas represent 1 values (high adhesive potency) and red areas represent 0 values (low adhesive potency).

5.3. Investigation of the hypothesis function

In section 3.6, we made the very informal argument that the model’s test error could be used to evaluate the valid-

ity of different hypothesis g describing the evolution of the copper surface characteristics with time t . In this section, we attempt to experimentally validate this idea.

We start by evaluating the linear function $g(t) = 1 - \frac{t}{T}$. This hypothesis function formulates the idea that for $t = 1$, all the segmentation mask values should be 1, so that the whole surface area should have strong adhesive potency, while for $t = 14$, all the segmentation mask values should be 0, so that the whole surface area should have low adhesive potency. For each t between 1 and 14, this hypothesis function defines a linear decrease of the low adhesive potency surface ratio. See Figure 5 for a visual illustration of the surface ratio evolution.

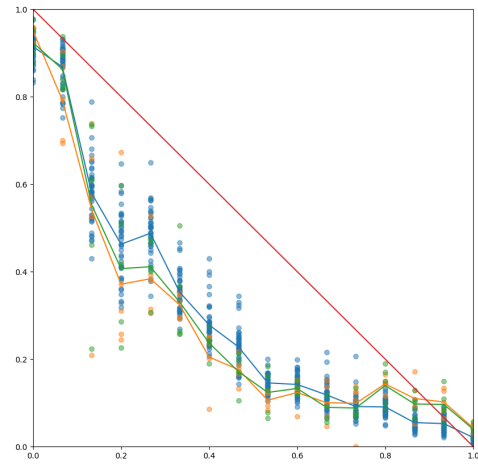


Figure 10. Results of the segmentation model evaluation. The red line illustrates the linear hypothesis used for training. Blue, yellow and green dots represent individual image output predictions on the training, validation and test set respectively. Even on training samples, the model is not able to overfit the ground labels generated by the hypothesis function. Instead of a linear decay, the model seems to predict a more aggressive decay, following a square root or logarithmic trend.

Figure 10 compares the hypothesis function to the actual model output on the training, validation and test set. First, we observe no overfitting as all datasets yield similar average results. Second, we observe that instead of following a linear trend from 1 to 0, the surface ratio seems to evolve more similarly to a squared root or logarithmic function with a sharp decrease in value for low t and a slower decrease in higher values of t . It is interesting to note that the same behavior is also observed on samples of the training set, for which the model was explicitly trained to reproduce the linear hypothesis function. This result suggests that the surface ratio (i.e. the smoothing effect of the stress test on the copper surface) may decrease sub-linearly with time.

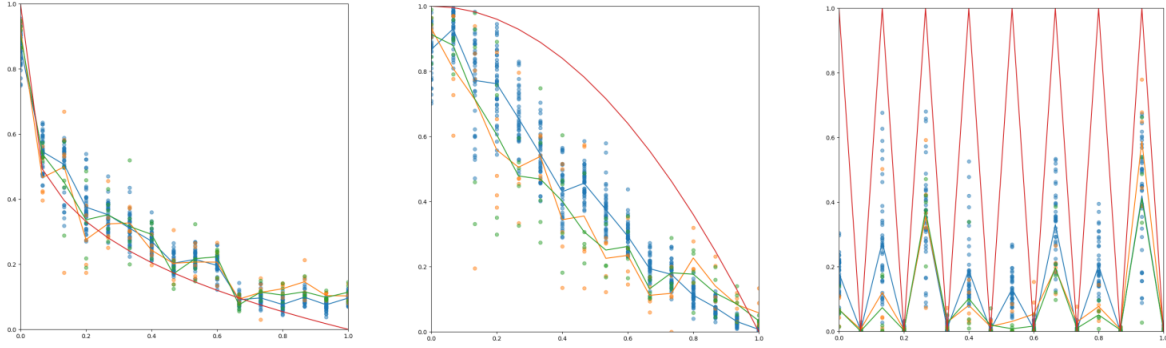


Figure 11. Results of the segmentation model evaluation with different hypothesis function g . (Left) $g = 1 - t^{-\frac{1}{2}}$. The model seems to nicely fit the hypothesis function. (Middle) $g = 1 - t^{-2}$. The model diverges from the hypothesis function, suggesting a more rapid decay. (Right) $g = \text{sign}(-1^t)$. The model struggles to learn the hypothesis function, yielding high errors.

Motivated by this first experiment, we train the model on different hypothesis functions g , illustrated in Figure 11. We find that training a model with a hypothesis functions g as suggested by our first experiment yield low model errors, while training the model with an unrealistic hypothesis (i.e. a sine function) yield high errors.

6. Discussion & Conclusion

In this paper, we have argued that deep neural networks can be used to help explain the process behind low-level recognition tasks. We have focused on the task of assessing the adhesive potency of copper surfaces in the context of PCB manufacturing. This is an interesting task to showcase the limitations of human’s ability to explain their visual process as it is a very low level recognition task for which trained experts can provide qualitative guesses while not being able to fully justify for their guess.

On this task, we have first shown that an unrestricted classifier architecture can outperform a human expert in accuracy. However, this architecture does not bring us any insight into *why* a given copper surface should be classified as a high or low adhesive potency surface.

To shed some light into the decision process of the model, we proposed to cast this problem as a a segmentation problem, in which the model is tasked with assigning a binary score to each pixel, indicating whether the local area represented by this pixel provides a high or low adhesive power. Visualizing the segmentation output of the model may prove useful for human experts to better identify surface patterns characteristic of high and low adhesive potency surfaces.

Finally, we developed a weak label training procedure to train this segmentation model. Our procedure relies on a hypothetical relationship between the manufacturing process and the copper surface roughness. We very informally argued that the model error may reflect the validity of the hypothesis function, and provided preliminary experiment

results that tend to confirm our informal argument: we observe lower errors when training the model with hypothesis functions that *seem* more plausible, and lower model error on unrealistic hypothesis functions (i.e. a sine function).

However, the result of this last experiment should be taken with caution as this study is still in a very preliminary stage. In particular, our training procedure relies on several assumptions that most likely do not hold in reality: We first assumed that each individual area could be represented by a single binary value defining its adhesive potency, and that the manufacturing process affects the mean value of the copper surface. In reality, the adhesive potency of local areas are most likely not binary in nature. Moreover, the adhesive potency may not be best measured as a spatial average over each local area.

Nevertheless, these preliminary results are encouraging, and we will continue our analysis in future work.

References

- [1] R. Ke, A. Bugeau, N. Papadakis, P. Schuetz, and C.-B. Schönlieb. A multi-task u-net for segmentation with lazy labels. *arXiv preprint arXiv:1906.12177*, 2019. 6
- [2] D. H. Park, L. A. Hendricks, Z. Akata, B. Schiele, T. Darrell, and M. Rohrbach. Attentive explanations: Justifying decisions and pointing to the evidence. *arXiv preprint arXiv:1612.04757*, 2016. 5
- [3] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 4
- [4] J. Yosinski, J. Clune, A. Nguyen, T. Fuchs, and H. Lipson. Understanding neural networks through deep visualization. *arXiv preprint arXiv:1506.06579*, 2015. 6