# Unsupervised Welding Defect Detection Using Audio and Video

Georg Stemmer[1], Jose A. Lopez[1*], Juan A. Del Hoyo Ontiveros[1], Arvind Raju[1], Tara Thimmanaik[1], Sovan Biswas[1]

[1*]Intel Corp., 2200 Mission College Blvd., Santa Clara, 95054, CA, USA.

*Corresponding author(s). E-mail(s): jose.a.lopez@intel.com;

**Abstract**

In this work we explore the application of AI to robotic welding. Robotic welding is a widely used technology in many industries, but robots currently do not have the capability to detect welding defects which get introduced due to various reasons in the welding process. We describe how deep-learning methods can be applied to detect weld defects in real-time by recording the welding process with microphones and a camera. Our findings are based on a large database with more than 4000 welding samples we collected which covers different weld types, materials and various defect categories. All deep learning models are trained in an unsupervised fashion because the space of possible defects is large and the defects in our data may contain biases. We demonstrate that a reliable real-time detection of most categories of weld defects is feasible both from audio and video, with improvements achieved by combining both modalities. Specifically, the multi-modal approach achieves an average Area-under-ROC-Curve (AUC) of 0.92 over all eleven defect types in our data. We conclude the paper with an analysis of the results by defect type and a discussion of future work.

## 1 Introduction

Robotic arc welding, i.e., the use of robots for automating the arc welding process, is a key manufacturing technology in many industries. As the quality of a weld depends on many factors, even a robot that repeats each step of the process perfectly will produce defective welds from time to time. The time at which a defect gets detected in the

welding process has direct impact on the cost associated with correcting the problem. So to minimize the overall defects and reduce the correction cost there is a growing interest in detecting anomalies in real-time, i.e., during the welding process, rather than post weld defect detection. Ideally, a future intelligent manufacturing system would be able to adjust the welding robot's parameters automatically even before the failure starts to occur in the welding process.

In this work, we investigate the use of deep learning models for automatic weld defect detection in real-time using a camera and microphones. Cameras can monitor the weld pool geometry and oscillations, which are known to be predictive of weld defects [1–3]. Cameras need to have line-of-sight to the weld pool, and this consequently limits the mounting options. Microphones, on the other hand, can capture audible disturbances in the welding process without line-of-sight. They have been shown to provide useful information about defects [4–10]. We investigate in the experimental section for several defect types whether they can be better detected visually or acoustically, and how the two modalities perform in combination.

Current deep learning models require large amounts of training data to estimate their parameters. We address this issue by collecting more than 4000 samples of good and defective welds, which is, to best of our knowledge, significantly more than what has been reported for similar datasets in the literature. Still, we consider the number of samples too small to train a classifier that can distinguish each relevant defect type reliably, given the variation of the input signals that can be expected in a real application. Therefore we decided to address the weld defect detection problem with an unsupervised anomaly detection approach. Anomaly detection models are trained on good welds only: all defects in our dataset occur solely during evaluation and are unseen in training. While this makes the defect specific performance more difficult to tune, we believe that the resulting performance is more representative of a real application.

Of course, we are not the first interested in bringing the benefits of AI to established spaces like manufacturing, and welding in particular. Companies have produced AI-based anomaly detection solutions in recent years [11–14], and numerous authors have explored the use of modern data-driven algorithms to improve weld defect detection [1, 15, 16], which can be costly to remediate [17]. In [18], the authors used sequence tagging and logistic regression to detect welding defects. Mohanasundari et al. [19] used post-weld images to classify defects. Buongiorno et al. [20] leveraged thermographic image sequences from an infrared camera to detect defects. Our work is also related to predictive maintenance solutions, which tend to use vibration or acoustic-emission sensors. Acoustic emission sensors have even been used successfully to monitor weld quality [18] as well. There is good reason for using this type of sensors for machine condition monitoring, as 90% of rotating machinery uses rolling-element bearings, which tend to be points of failure [21]. On the other hand, welding is a much more complex process. Therefore we expect that microphones and cameras are better suited for detecting weld defects than acoustic emission or vibration sensors.

The main contributions of this work are the following: We collect a large multi-modal dataset of samples of robotic arc welding in a real industrial environment. It covers different weld types, various welding parameter configurations, and steel types.

The size of the dataset allows us to train deep ML models, as opposed to shallow or analytical models often described in the literature (e.g., [18]). With these deep ML models we are able to demonstrate that camera and microphone are adequate sensors for real-time weld defect detection. Our unsupervised approach ensures that the model is not biased to defect characteristics which are specific for our data collection setup and allows us to compare results between different defect types and and modalities. Finally, we demonstrate that defect detection performance can be improved using a multi-modal combination of both sensor types.

In the next sections we provide a comprehensive description of how the dataset has been collected. We introduce our experimental setup, quality metrics and ML models. The performance of unsupervised algorithms trained in single and multi-modal fashion on this dataset is evaluated experimentally and we present our conclusions.

## 2 Data Set

The goal of our data collection was to record enough samples for each of the most important weld categories to make statistically valid comparisons across modalities and weld defect types. We collaborated with a supplier that has access to automotive factories to conduct the data acquisition in a real factory environment. The welds were generated using a 6-axis arc welding robot using two steel types and thicknesses. The steel types were selected to be often-used varieties for automotive applications in India. The first type is known as "FE410" and the second "BSK46" type has higher carbon content and is used for higher-strength applications. The thicknesses used were 7mm for most of the samples and 3mm for specific defects that could not be efficiently induced using the 7mm steel, like burnthrough.

### 2.1 Recording Setup, Data Collection Procedure, and Limitations

The data collection station comprised an AII-V6 6-axis arc welding robot [28], an Intel i7-based workstation and camera [29], two high-bandwidth microphones [25], and an audio interface [23]. The microphones were selected to enable studies on the observability of weld defects at higher frequencies. Moreover, audio samples were recorded using a 192 KHz sampling rate and saved in lossless FLAC format [26]. The video samples were recorded at a nominal 30 FPS and saved in AVI format. Fig. 1 shows the KML camera and computer. The camera was attached to the welder arm, about 200mm from the torch. The microphones were attached to the work table, about 300mm from the welder arm motion axis.

To generate weld defects, the supplier contracted a welding expert who supervised the initial configuration of the welding robot to produce the desired weld defects. The welding expert did not perform any post-weld validation or labeling. Thus, the dataset contains some amount of label noise. Moreover, all sensors were triggered to start recording at the same time, but there is some variation in the actual time it takes for the individual device to respond. To align the different modalities which is required for multi-modal experiments, the weld start and end times were identified for each modality by inspecting the audio and illumination changes at the start and end

(a) KML camera.



(b) KML control panel.



(c) KML computer.

**Fig. 1**: KML camera and computer.

of welding. The audio data, which was collected through a dedicated audio interface with a low-latency driver [23] was taken as the ground truth source for determining the welding duration. In this way, it was determined that the actual recorded frames-per-second (FPS) of the video varies from the expected 30 FPS. It is worth mentioning that this post-collection alignment would not be needed if the robot command signals were readily accessible.

| Quantity | Component | Description |
|---|---|---|
| 1 | welding robot | OTC AII-V6 6-axis arc welding robot |
| 1 | camera | KML welding camera |
| 2 | microphone | Earthworks SR314 |
| 1 | microphone interface | MOTU M4 audio interface |
| 1 | workstation | KML workstation with Intel i7 processor |

**Table 1**: Data collection station parts list.

## 2.2 Welding Sample Distribution

The dataset contains weld samples for the 12 weld categories shown in Tab. 2. Each weld category has been recorded for different weld types, and, where applicable, for different materials. Tab. 3 contains a brief summary of the dataset by weld type. A complete breakdown of the dataset is included in Appendix A. Fig. B1 in Appendix B shows some examples of (post-weld) photos of welding samples.

| Weld Category | Total |
|---|---|
| Good | 819 |
| Excessive Convexity | 160 |
| Undercut | 160 |
| Crater Cracks | 161 |
| Overlap | 160 |
| Excessive Penetration | 480 |
| Porosity w/Excessive Penetration | 480 |
| Spatter | 320 |
| Lack Of Fusion | 320 |
| Warping | 320 |
| Porosity | 340 |
| Burnthrough | 320 |
| **All** | 4040 |

**Table 2**: Summary of the number of weld samples by weld category.

| Weld Type | Samples | Material | Samples | Total |
|---|---|---|---|---|
| non-fillet | 1720 | 7mm-FE410 | 2560 | 2919 |
| non-fillet | 840 | 3mm-FE410 | | |
| non-fillet | 239 | 7mm-BSK46 | 359 | |
| non-fillet | 120 | 3mm-BSK46 | | |
| fillet | 701 | 7mm-FE410 | 981 | 1121 |
| fillet | 280 | 3mm-FE410 | | |
| fillet | 100 | 7mm-BSK46 | 140 | |
| fillet | 40 | 3mm-BSK46 | | |

**Table 3**: Summary of the number of weld samples by weld type and material used.

# 3 Experimental Setup

The experiments in this work focus on the following questions: (1) Can weld defects be reliably detected in real-time from audio and/or video recordings? (2) How does the detection accuracy depend on the defect type and modality? (3) What improvements can be expected from combining audio and video in a multi-modal system?

To answer those questions, we treat the weld defect problem as an anomaly detection problem, not a classification problem. That means, our models do not attempt to classify a sample into a weld category, but generate an anomaly score that increases with the likelihood of a defect. For a real use case a threshold has to be determined in advance. If the score of the model exceeds this threshold during welding, a defect will be detected and the supervisor of the robot will be notified. The choice of the threshold depends on the false positive and false negative rates of the model for the defect types that are relevant for a specific use case, the costs associated with a false detection of a defect, and the costs resulting from a missed defect. As we would like to evaluate the quality of our models independently of an application-specific threshold, we compare the models using the Area-under-ROC-Curve (AUC) metric which is scale- and threshold-independent [39]. AUC is defined as the area under the Receiver Operating Characteristic (ROC) curve for all possible values of the false positive rate (FPR):

$$AUC = \int_{x=0}^{1} ROC(x)dx \qquad (1)$$

where $x$ is the FPR and $ROC(x)$ is the true positive rate (TPR) [24].

When calculating FPR and TPR for a specific threshold on the validation and test sets, we have to take into account that we are targeting a real-time scenario, where the model produces an anomaly score at every time frame. On the other hand, each sample in our data set has just a single label denoting the weld category. For instance, a sample could be labeled as "porosity" but there is no indication at which time instance the defect occurs. Therefore, we aggregate all scores produced by the model for all frames of the sample by either taking their maximum or expected value.

In this work we employ simple but effective convolutional neural networks (CNN) for weld defect detection. This allows us to focus more on gaining insights and less on conducting large hyper-parameter optimizations and architecture searches. For audio, we used a 1D CNN auto-encoder [31]. For video, we applied the relatively simple 3D CNN from [35] provided by the MMAction2 library [36][1]. All our work used the PyTorch tensor library for Python [32, 34].

For the experiments, the dataset is split into a training, validation, and test partition. The training data includes only good, i.e., normal-state, samples, while the defects are divided equally among the validation and test subsets. The validation partition of the dataset is used for hyperparameter tuning. The best hyperparameter configuration on the validation partition is used to generate the test results from the test partition. The data split is summarized in Tab. 4. The same partitions were used for all experiments with both modalities.

---

[1]The "slowfast_r101_4x16x1_256e_kinetics400_rgb_20210218-d8b58813.pth" checkpoint from v0.15.0 was used [30].

| Partition | Number of "Good" Samples | Number of "Defective" Samples |
|---|---|---|
| Train | 576 | 0 |
| Validation | 122 | 1610 |
| Test | 121 | 1611 |

**Table 4**: Data split into training, validation, and test partitions.

## 3.1 Acoustic Anomaly Detection

For the audio experiments, the original sampling rate of 192 kHz was maintained and only one channel was used. The audio CNN auto-encoder architecture is shown in Table 5. A detailed description of the model's topology can be found in [31]: the key characteristic is that the bottleneck layer largely preserves the time dimension. This feature has been motivated by the work of Agrawal et al. [33]. The model in Tab. 5 is much smaller than the one described in [33] because the convolutions are not gated.

| Layer | Input, Output Channels | Kernel Size | Stride |
|---|---|---|---|
| BatchNorm1D | n-bins | N/A | N/A |
| Conv1D | (n-bins, 1024) | 3 | 1 |
| 3 x Conv1D | (1024, 1024) | 3 | 1 |
| Conv1D | (1024, bottleneck size) | 3 | 1 |
| ConvTranspose1D | (bottleneck size, 1024) | 3 | 1 |
| 3 x ConvTranspose1D | (1024, 1024) | 3 | 1 |
| ConvTranspose1D | (1024, n-bins) | 3 | 1 |

**Table 5**: Audio 1D CNN.

The model uses leaky-ReLU activations in all but the last activation before the output layer, which uses a PReLU activation. Overall, the model has 31,670,306 trainable parameters.

The latency of this model equals the hop length used to train the model. For example, for a hop length of 8192 and 192 kHz audio, the latency is about 42.7 ms. The model has 5 encoding layers with kernel size 3. Each encoding layer decreases the time dimension by 2, therefore inputs must have more than 10 frames. This means the input buffer must satisfy Eq. 2, with the FFT window an integer multiple of the hop length. In the foregoing example, the input buffer needs to hold $12 \times 8192$ samples.

$$\text{buffer size} = \text{hop length} \times \left( 10 + \frac{\text{FFT window}}{\text{hop length}} \right) \quad (2)$$

## 3.2 Visual Anomaly Detection

Our approach to visual weld defect detection is based on a two-stage process. The first stage encodes each video frame into a fixed dimensional feature vector. For this, a window that is 64 frames long is shifted frame-by-frame over the whole video. 64 frames correspond to roughly 2 seconds at a frame rate of 30 FPS. As the window is centered around the frame of interest, the 64-frame window size leads to defect detection latency of around one second in a real-time scenario. This is a larger latency

than for the acoustic anomaly detection, which is based on much smaller windows as will be described in Sec. 4. Still, we believe that it should be acceptable for many use cases. Each window is encoded into a 2304-dimensional feature vector using the pre-trained and fixed Slowfast [35] model. In the second stage, an auto-encoder model consisting of an encoder, a bottleneck, and a decoder is used for generating anomaly scores from the input feature vector. The encoder maps the 2304-dimensional input feature vector to 64-dimensional latent space by passing it through multiple linear layers along with ReLU activation and dropout in sequence. Each linear layer reduces the dimension by $\frac{1}{2}$, thus creating a bottleneck with a 64-dimensional latent vector. Later, the decoder uses the latent space embedding to reconstruct the original feature of the frame. The decoder consists of multiple linear layers with ReLU activation and dropout as well. Each linear layer of the decoder scales the dimension by 2 until the original 2304 dimension is obtained. All dropout layers zero weights with probability 0.5. The architecture details of the auto-encoder are provided in Tab. 6.

The auto-encoder model is trained on welding videos minimizing the model's anomaly score which is defined as the mean-squared-error (MSE) between the input feature vector and the output of the decoder. Note, that the Slowfast model from the first stage is fixed, and no-back propagation is applied to the first stage during training.

| Layer | Input, Output Dim. |
|---|---|
| Linear | (2304, 512) |
| Linear | (512, 256) |
| Dropout ($p = 0.5$) | - |
| Linear | (256, 128) |
| Dropout ($p = 0.5$) | - |
| Linear | (128, 64) |
| Dropout ($p = 0.5$) | - |
| Linear | (64, 64) |
| Linear | (64, 64) |
| Linear | (64, 128) |
| Linear | (128,256) |
| Linear | (256, 512) |
| Dropout ($p = 0.5$) | - |
| Linear | (512, 2304) |

**Table 6**: Auto-encoder model for the video modality.

## 3.3 Multi-modal Anomaly Detection

We used a late-fusion approach to combine the anomaly scores of the two modalities. Since anomaly scores generated by different models generally have different scales, we first standardized the scores using the mean and standard deviation of the anomaly scores over the training set. Next, we identified the optimal convex combination of the audio and video anomaly scores by running grid search over convex combinations on the validation data. The best weighting is applied to compute the anomaly scores on the test set and to produce the final results.

# 4 Results

## 4.1 Acoustic Anomaly Detection

For hyperparameter tuning we trained separate auto-encoder models with different parameter settings on the training partition and evaluated their AUC on the validation set. More specifically, a grid search was performed over FFT window sizes 4096, 16384, 32768, and 65536 and bottleneck dimensions 16, 32, 48, and 64. The hop length was fixed at 50% of the FFT window size. We remind the reader that for a given analysis (time) window, the corresponding FFT window length is proportional to the sampling rate. Therefore, the FFT window sizes included in the grid search range from approximately 21 ms to 341 ms for the sampling rate in our data set, which is 192 kHz. The latency of these models range from approximately 11ms to 171ms. All models were trained for 50 epochs using a one-cycle learning schedule with a peak learning rate of $1 \times 10^{-4}$ for the Adam optimizer [37, 38]. The MSE loss was minimized during training.

Tab. 7 shows the FFT window search results. These experiments determine that an FFT window of 16384 with a bottleneck dimension of 48 produced the best performance on the validation set. To obtain the AUC scores we explored several frame-wise anomaly score aggregation methods: expected value, moving average (MA) smoothing, and taking the max. We found taking the average worked best, followed closely by MA smoothing, and taking taking the max was last. We only show the expected value scores for space considerations.

| FFT Window | Val AUC (b.n.=16) | Val AUC (b.n.=32) | Val AUC (b.n.=48) | Val AUC (b.n.=64) |
|---|---|---|---|---|
| 4096 | 0.7550 | 0.7953 | 0.7597 | 0.7549 |
| 16384 | **0.8444** | **0.8434** | **0.8451** | **0.8426** |
| 32768 | 0.8210 | 0.8231 | 0.8262 | 0.8187 |
| 65536 | 0.8193 | 0.8206 | 0.8215 | 0.8234 |

**Table 7**: FFT Window and Bottleneck Search.

On the test set, the best validation model obtained an AUC of 0.8460. This is not much different from the AUC on the validation set which indicates that the model does not overfit. To gain insights into the difficulty of detecting particular defects, we also determined the validation and test AUC by defect type – the results are shown in Tab. 8 and Tab. 9, respectively. Again, the results indicate a high agreement between validation and test set AUCs. Note that these results are not directly actionable because in most use cases the defect type is not known a priori. Breaking-out performance by defect type is, however, useful for estimating performance gains that could be obtained when fusing audio with video information.

In cases where an error (or error cost) distribution is known a priori, one may select the optimal hyperparameters differently. Additionally, latency or memory requirements can affect the FFT window and hop length selection as well.

| Welding Category | FFT=4096 | FFT=16384 | FFT=32768 | FFT=65536 |
|---|---|---|---|---|
| Excessive Penetration | 0.7023 | 0.7717 | 0.8157 | **0.8268** |
| Burnthrough | **0.8256** | 0.6920 | 0.5952 | 0.6256 |
| Porosity | **0.9898** | 0.9892 | 0.9633 | 0.9385 |
| Porosity w/EP | 0.9584 | **0.9758** | 0.9501 | 0.9460 |
| Undercut | **0.9527** | 0.9143 | 0.8957 | 0.9000 |
| Overlap | 0.9800 | 0.9840 | 0.9863 | **0.9874** |
| Lack of fusion | 0.4702 | 0.8316 | **0.8345** | 0.8226 |
| Excessive Convexity | **0.8981** | 0.8845 | 0.8704 | 0.8548 |
| Spatter | 0.3591 | 0.7773 | 0.8531 | **0.8906** |
| Warping | 0.6411 | **0.7518** | 0.5925 | 0.5030 |
| Crater Cracks | **0.9912** | 0.9672 | 0.9678 | 0.9659 |
| **All** | 0.7597 | **0.8451** | 0.8262 | 0.8215 |

**Table 8**: Validation AUC by FFT window size and defect type for model with bottleneck=48.

| Welding Category | FFT=4096 | FFT=16384 | FFT=32768 | FFT=65536 |
|---|---|---|---|---|
| Excessive Penetration | 0.7023 | 0.7807 | **0.8220** | 0.8204 |
| Burnthrough | **0.8250** | 0.7018 | 0.6202 | 0.6497 |
| Porosity | 0.9896 | **0.9949** | 0.956 | 0.9144 |
| Porosity w/EP | **0.9663** | 0.9659 | 0.9360 | 0.9155 |
| Undercut | **0.9489** | 0.8989 | 0.8811 | 0.8637 |
| Overlap | 0.9695 | **0.9743** | 0.9671 | 0.9631 |
| Lack of fusion | 0.4693 | **0.8372** | 0.8353 | 0.8211 |
| Excessive Convexity | **0.8921** | 0.8842 | 0.8553 | 0.8361 |
| Spatter | 0.3855 | 0.7796 | 0.8407 | **0.8678** |
| Warping | 0.6483 | **0.7648** | 0.6723 | 0.6060 |
| Crater Cracks | **0.9935** | 0.9453 | 0.9218 | 0.9093 |
| **All** | 0.7651 | **0.8460** | 0.8296 | 0.8157 |

**Table 9**: Test AUC by FFT window size and defect type for model with bottleneck=48.

## 4.2 Visual Anomaly Detection

For visual defect detection, we first generate feature vectors of dimension 2304 using the pre-trained Slowfast model for each sample in all data partitions. In a second step we train the auto-encoder model using the Adam optimizer with a learning rate of 0.0005 and MSE loss function for up to 1000 epochs on the feature vectors. For the model that performs best on the validation set, we compare different methods to aggregate the scores: simply taking the maximum score ("Max w/o smoothing"), smoothing the scores by averaging them within a one-second window ("Max over 1s-MA"), or smoothing the scores by averaging them within the full two-second window ("Max over 2s-MA"). Table 11 shows the AUC for different aggregation methods on the test set. The best validation model obtains an overall AUC of 0.9052 and 0.8977 on the validation and test data, respectively.

| Weld Category | Max w/o Smoothing (AUC) | Max Over 1s-MA (AUC) | Max Over 2s-MA (AUC) |
|---|---|---|---|
| Excessive Penetration | 0.8165 | 0.8815 | 0.9011 |
| Burnthrough | 0.8066 | 0.8978 | 0.9156 |
| Porosity | 0.9732 | 0.9984 | 0.9998 |
| Porosity w/EP | 0.9701 | 0.9985 | 0.9999 |
| Undercut | 0.9180 | 0.9462 | 0.9575 |
| Overlap | 0.7446 | 0.7973 | 0.8155 |
| Lack of fusion | 0.7229 | 0.7461 | 0.7671 |
| Excessive Convexity | 0.9809 | 0.9965 | 0.9984 |
| Spatter | 0.9458 | 0.9793 | 0.9857 |
| Warping | 0.6690 | 0.6991 | 0.7141 |
| Crater Cracks | 0.9020 | 0.8597 | 0.8635 |
| **All** | 0.8577 | 0.8947 | **0.9052** |

**Table 10**: Defect specific AUC for various anomaly score aggregation methods calculated for the validation data using video.

| Weld Category | Max w/o Smoothing (AUC) | Max Over 1s-MA (AUC) | Max Over 2s-MA (AUC) |
|---|---|---|---|
| Excessive Penetration | 0.7746 | 0.8534 | 0.8747 |
| Burnthrough | 0.7674 | 0.8841 | 0.8997 |
| Porosity | 0.9557 | 0.9847 | 0.9868 |
| Porosity w/EP | 0.9558 | 0.9846 | 0.9873 |
| Undercut | 0.8744 | 0.9159 | 0.9270 |
| Overlap | 0.7530 | 0.8319 | 0.8516 |
| Lack of fusion | 0.7055 | 0.7644 | 0.7784 |
| Excessive Convexity | 0.9673 | 0.9825 | 0.9844 |
| Spatter | 0.9265 | 0.9583 | 0.9654 |
| Warping | 0.6632 | 0.7375 | 0.7456 |
| Crater Cracks | 0.8527 | 0.8323 | 0.8436 |
| **All** | 0.8345 | 0.8873 | **0.8977** |

**Table 11**: Defect specific AUC for various anomaly score aggregation methods calculated for the test data using video.

## 4.3 Multi-modal Anomaly Detection

The input scores for the late fusion combination of the video and audio modalities were created using the best performing models for each modality, i.e., the audio model using FFT window size 16384 and bottleneck dimension of 48, and the video model with 2s-MA smoothing. We determined the weighting for the modalities on the validation data using the grid search described in Sec. 3.3 with step size 0.01.

This way, we found the best weighting to be 0.37 and 0.63 for the audio and video scores, respectively. These weightings make sense given the stronger performance of the video modality.

Using the optimal weighting we obtained our overall best test AUC of 0.9178. From Tab. 12, it can be seen that the overall audio validation and test AUCs improved by 9.8% and 8.5%, respectively, and the overall video validation and test AUCs improved by 2.5% and 2.2%, respectively. The defect-specific scores improved as well. For audio,

the AUC metric improved for 8 of 11 defect categories. For video, the AUC improved for 7 of 11 categories. It is worth mentioning that the weakest performing defect categories all improved. In particular, audio had four defect categories with AUCs in the 0.7s and all benefited from fusing. For video, two defect categories had AUCs in the 0.7s and both benefited from fusing. See the entries of Tab. 13 shown in bold.

Fig. 2 shows the detection error curves for the multi-modal predictions, on the test data. As can be seen, the FPR and false-negative-rate (FNR) intersect at about 17%.

| Weld Category | Validation (AUC) | Test (AUC) |
|---|---|---|
| Excessive Penetration | 0.9134 | 0.8893 |
| Burnthrough | 0.8811 | 0.8713 |
| Porosity | 1 | 0.9901 |
| Porosity w/EP | 0.9997 | 0.9884 |
| Undercut | 0.9800 | 0.9477 |
| Overlap | 0.9573 | 0.9310 |
| Lack of fusion | 0.8306 | 0.8455 |
| Excessive Convexity | 0.9982 | 0.9823 |
| Spatter | 0.9809 | 0.9576 |
| Warping | 0.7671 | 0.7984 |
| Crater Cracks | 0.9567 | 0.9266 |
| **All** | 0.9280 | 0.9178 |

**Table 12**: Multi-modal performance on validation and test data.

| Weld Category | Audio Change % | Video Change % |
|---|---|---|
| Excessive Penetration | **13.9215** | 1.6732 |
| Burnthrough | **24.1646** | -3.1575 |
| Porosity | 3.9451 | 0.3350 |
| Porosity w/EP | 5.7199 | 0.1151 |
| Undercut | 5.5453 | 2.2401 |
| Overlap | -4.2601 | 9.3291 |
| Lack of fusion | 0.9933 | **8.6198** |
| Excessive Convexity | 11.1384 | -0.2099 |
| Spatter | **22.8806** | -0.8132 |
| Warping | **4.9216** | **7.0805** |
| Crater Cracks | -1.6248 | 9.8452 |
| **All** | 8.4811 | 2.2337 |

**Table 13**: Multi-modal test data performance difference by modality and defect. The bolded entries show the changes on defects with unimodal AUCs in the 0.7s.

# 5 Conclusions and Future Work

In this work, we explored unsupervised weld defect detection using audio, video, and their combination. We demonstrated that, using a deep-learning based approach, both

**Fig. 2**: Accept-reject probability of multi-modal predictions on test data.

modalities allow to reliably detect the most important defect types in real-time. The best audio model had a latency of 42.7ms and is best suited for detecting porosity, overlap, and crater cracks. It has the lowest AUC scores for burnthrough, excessive penetration, spatter, and warping. Video generally shows better AUC than audio, but this comes at the price of a larger overall model size and a higher latency. Using video, the best detection performance is achieved for porosity, excessive convexity, and spatter. We observe the lowest performance for lack of fusion, and warping. We demonstrated that a combined approach using late fusion of normalized scores for both modalities offers improvements. More specifically, AUC scores for all defect types average 0.92, with the lowest score approximately 0.80. The worst performing categories for both modalities all improved.

Future work will investigate more elaborate ways to combine the two modalities, e.g., a joint model that directly incorporates input from all sensors. Furthermore, we believe that the biggest limitation of our dataset is that it has been collected in a supervised way, i.e., the robot has been configured by purpose to generate defects as we required a sufficient number of defect samples for our experiments. This approach creates a potential mismatch to a real use case where defects occur randomly and rarely. We expect that defects will show up less pronounced than in our dataset. Therefore, we plan to record a real production process of a collaboration partner and to label the defects as they occur. This will allow us to validate our proposed methods in a more realistic setting.

## Compliance with Ethical Standards

# Appendix A   Sample Distributions

This appendix contains a complete dataset description by defect category: Tab. A1 shows the distribution of the samples for the normal-state, i.e., good weld category, by weld type and material. Tab. A2 to A12 show the sample distribution for the remaining weld categories.

| Weld Category | Weld Type | Samples | Material | Samples | Total |
|---|---|---|---|---|---|
| Good | non-fillet | 580 | 7mm-FE410 | 580 | 659 |
| Good | non-fillet | 0 | 3mm-FE410 | | |
| Good | non-fillet | 79 | 7mm-BSK46 | 79 | |
| Good | non-fillet | 0 | 3mm-BSK46 | | |
| Good | fillet | 140 | 7mm-FE410 | 140 | 160 |
| Good | fillet | 0 | 3mm-FE410 | | |
| Good | fillet | 20 | 7mm-BSK46 | 20 | |
| Good | fillet | 0 | 3mm-BSK46 | | |

**Table A1**: The normal-state welding sample distribution.

| Weld Category | Weld Type | Samples | Material | Samples | Total |
|---|---|---|---|---|---|
| Excessive Convexity | non-fillet | 0 | 7mm-FE410 | 0 | 0 |
| Excessive Convexity | non-fillet | 0 | 3mm-FE410 | | |
| Excessive Convexity | non-fillet | 0 | 7mm-BSK46 | 0 | |
| Excessive Convexity | non-fillet | 0 | 3mm-BSK46 | | |
| Excessive Convexity | fillet | 140 | 7mm-FE410 | 140 | 160 |
| Excessive Convexity | fillet | 0 | 3mm-FE410 | | |
| Excessive Convexity | fillet | 20 | 7mm-BSK46 | 20 | |
| Excessive Convexity | fillet | 0 | 3mm-BSK46 | | |

**Table A2**: The distribution of welding samples containing excessive convexity.

| Weld Category | Weld Type | Samples | Material | Samples | Total |
|---|---|---|---|---|---|
| Undercut | non-fillet | 0 | 7mm-FE410 | 0 | 0 |
| Undercut | non-fillet | 0 | 3mm-FE410 | | |
| Undercut | non-fillet | 0 | 7mm-BSK46 | 0 | |
| Undercut | non-fillet | 0 | 3mm-BSK46 | | |
| Undercut | fillet | 140 | 7mm-FE410 | 140 | 160 |
| Undercut | fillet | 0 | 3mm-FE410 | | |
| Undercut | fillet | 20 | 7mm-BSK46 | 20 | |
| Undercut | fillet | 0 | 3mm-BSK46 | | |

**Table A3**: The distribution of welding samples containing undercut.

| Weld Category | Weld Type | Samples | Material | Samples | Total |
|---|---|---|---|---|---|
| Crater Cracks | non-fillet | 0 | 7mm-FE410 | 0 | 0 |
| Crater Cracks | non-fillet | 0 | 3mm-FE410 | | |
| Crater Cracks | non-fillet | 0 | 7mm-BSK46 | 0 | |
| Crater Cracks | non-fillet | 0 | 3mm-BSK46 | | |
| Crater Cracks | fillet | 141 | 7mm-FE410 | 141 | 161 |
| Crater Cracks | fillet | 0 | 3mm-FE410 | | |
| Crater Cracks | fillet | 20 | 7mm-BSK46 | 20 | |
| Crater Cracks | fillet | 0 | 3mm-BSK46 | | |

**Table A4**: The distribution of welding samples containing crater cracks.

| Weld Category | Weld Type | Samples | Material | Samples | Total |
|---|---|---|---|---|---|
| Overlap | non-fillet | 0 | 7mm-FE410 | 0 | 0 |
| Overlap | non-fillet | 0 | 3mm-FE410 | | |
| Overlap | non-fillet | 0 | 7mm-BSK46 | 0 | |
| Overlap | non-fillet | 0 | 3mm-BSK46 | | |
| Overlap | fillet | 140 | 7mm-FE410 | 140 | 160 |
| Overlap | fillet | 0 | 3mm-FE410 | | |
| Overlap | fillet | 20 | 7mm-BSK46 | 20 | |
| Overlap | fillet | 0 | 3mm-BSK46 | | |

**Table A5**: The distribution of welding samples containing overlap.

| Weld Category | Weld Type | Samples | Material | Samples | Total |
|---|---|---|---|---|---|
| Excessive Penetration | non-fillet | 0 | 7mm-FE410 | 280 | 320 |
| Excessive Penetration | non-fillet | 280 | 3mm-FE410 | | |
| Excessive Penetration | non-fillet | 0 | 7mm-BSK46 | 40 | |
| Excessive Penetration | non-fillet | 40 | 3mm-BSK46 | | |
| Excessive Penetration | fillet | 0 | 7mm-FE410 | 140 | 160 |
| Excessive Penetration | fillet | 140 | 3mm-FE410 | | |
| Excessive Penetration | fillet | 0 | 7mm-BSK46 | 20 | |
| Excessive Penetration | fillet | 20 | 3mm-BSK46 | | |

**Table A6**: The distribution of welding samples containing excessive penetration.

# Appendix B    Sample Images

This Appendix shows in Fig. B1 examples of photos taken of samples after the welding has been completed.

# References

[1] N. Lv, Y. Xu, S. Li, X. Yu, S. Chen, Deep learning based real-time and in-situ monitoring of weld penetration: Where we are and what are needed revolutionary solutions? Journal of Manufacturing Processes **93**, 15–46 (2023). https://doi.org/https://doi.org/10.1016/j.jmapro.2023.03.011

| Weld Category | Weld Type | Samples | Material | Samples | Total |
|---|---|---|---|---|---|
| Porosity w/ Excessive Penetration | non-fillet | 0 | 7mm-FE410 | 280 | 320 |
| Porosity w/ Excessive Penetration | non-fillet | 280 | 3mm-FE410 | | |
| Porosity w/ Excessive Penetration | non-fillet | 0 | 7mm-BSK46 | 40 | |
| Porosity w/ Excessive Penetration | non-fillet | 40 | 3mm-BSK46 | | |
| Porosity w/ Excessive Penetration | fillet | 0 | 7mm-FE410 | 140 | 160 |
| Porosity w/ Excessive Penetration | fillet | 140 | 3mm-FE410 | | |
| Porosity w/ Excessive Penetration | fillet | 0 | 7mm-BSK46 | 20 | |
| Porosity w/ Excessive Penetration | fillet | 20 | 3mm-BSK46 | | |

**Table A7**: The distribution of welding samples containing porosity and excessive penetration.

| Weld Category | Weld Type | Samples | Material | Samples | Total |
|---|---|---|---|---|---|
| Spatter | non-fillet | 280 | 7mm-FE410 | 280 | 320 |
| Spatter | non-fillet | 0 | 3mm-FE410 | | |
| Spatter | non-fillet | 40 | 7mm-BSK46 | 40 | |
| Spatter | non-fillet | 0 | 3mm-BSK46 | | |
| Spatter | fillet | 0 | 7mm-FE410 | 0 | 0 |
| Spatter | fillet | 0 | 3mm-FE410 | | |
| Spatter | fillet | 0 | 7mm-BSK46 | 0 | |
| Spatter | fillet | 0 | 3mm-BSK46 | | |

**Table A8**: The distribution of welding samples containing spatter.

| Weld Category | Weld Type | Samples | Material | Samples | Total |
|---|---|---|---|---|---|
| Lack Of Fusion | non-fillet | 280 | 7mm-FE410 | 280 | 320 |
| Lack Of Fusion | non-fillet | 0 | 3mm-FE410 | | |
| Lack Of Fusion | non-fillet | 40 | 7mm-BSK46 | 40 | |
| Lack Of Fusion | non-fillet | 0 | 3mm-BSK46 | | |
| Lack Of Fusion | fillet | 0 | 7mm-FE410 | 0 | 0 |
| Lack Of Fusion | fillet | 0 | 3mm-FE410 | | |
| Lack Of Fusion | fillet | 0 | 7mm-BSK46 | 0 | |
| Lack Of Fusion | fillet | 0 | 3mm-BSK46 | | |

**Table A9**: The distribution of welding samples with insufficient fusion.

[2] G. Ma, H. Yuan, L. Yu, Y. He, Monitoring of weld defects of visual sensing assisted gmaw process with galvanized steel. Materials and Manufacturing Processes **36**(10), 1178–1188 (2021). https://doi.org/10.1080/10426914.2021.1885711

[3] S. Zou, Z. Wang, S. Hu, W. Wang, Y. Cao, Control of weld penetration depth using relative fluctuation coefficient as feedback. Journal of Intelligent Manufacturing **31**(5), 1203–1213 (2020). https://doi.org/10.1007/s10845-019-01506-8

[4] T. Ji, N. Mohamad Nor, Deep learning-empowered digital twin using acoustic signal for welding quality inspection. Sensors **23**(5) (2023). https://doi.org/10.3390/s23052643

| Weld Category | Weld Type | Samples | Material | Samples | Total |
|---|---|---|---|---|---|
| Warping | non-fillet | 280 | 7mm-FE410 | 280 | 320 |
| Warping | non-fillet | 0 | 3mm-FE410 | | |
| Warping | non-fillet | 40 | 7mm-BSK46 | 40 | |
| Warping | non-fillet | 0 | 3mm-BSK46 | | |
| Warping | fillet | 0 | 7mm-FE410 | 0 | 0 |
| Warping | fillet | 0 | 3mm-FE410 | | |
| Warping | fillet | 0 | 7mm-BSK46 | 0 | |
| Warping | fillet | 0 | 3mm-BSK46 | | |

**Table A10**: The distribution of welding samples containing warping.

| Weld Category | Weld Type | Samples | Material | Samples | Total |
|---|---|---|---|---|---|
| Porosity | non-fillet | 300 | 7mm-FE410 | 300 | 340 |
| Porosity | non-fillet | 0 | 3mm-FE410 | | |
| Porosity | non-fillet | 40 | 7mm-BSK46 | 40 | |
| Porosity | non-fillet | 0 | 3mm-BSK46 | | |
| Porosity | fillet | 0 | 7mm-FE410 | 0 | 0 |
| Porosity | fillet | 0 | 3mm-FE410 | | |
| Porosity | fillet | 0 | 7mm-BSK46 | 0 | |
| Porosity | fillet | 0 | 3mm-BSK46 | | |

**Table A11**: The distribution of welding samples containing porosity.

| Weld Category | Weld Type | Samples | Material | Samples | Total |
|---|---|---|---|---|---|
| Burnthrough | non-fillet | 0 | 7mm-FE410 | 280 | 320 |
| Burnthrough | non-fillet | 280 | 3mm-FE410 | | |
| Burnthrough | non-fillet | 0 | 7mm-BSK46 | 40 | |
| Burnthrough | non-fillet | 40 | 3mm-BSK46 | | |
| Burnthrough | fillet | 0 | 7mm-FE410 | 0 | 0 |
| Burnthrough | fillet | 0 | 3mm-FE410 | | |
| Burnthrough | fillet | 0 | 7mm-BSK46 | 0 | |
| Burnthrough | fillet | 0 | 3mm-BSK46 | | |

**Table A12**: The distribution of welding samples containing burnthrough.

[5] J.Y.I. Alcaraz, W. Foqué, A. Sharma, T. Tjahjowidodo, Indirect porosity detection and root-cause identification in waam. Journal of Intelligent Manufacturing (2023). https://doi.org/10.1007/s10845-023-02128-x

[6] L. Chen, X. Yao, C. Tan, W. He, J. Su, F. Weng, Y. Chew, N.P.H. Ng, S.K. Moon, In-situ crack and keyhole pore detection in laser directed energy deposition through acoustic signal and deep learning. Additive Manufacturing **69**, 103547 (2023). https://doi.org/10.1016/j.addma.2023.103547

[7] A. Madhvacharyula, A. Pavan, S. Gorthi, S. Chitral, N. Venkaiah, D. Kiran, In situ detection of welding defects: a review. Welding in the World **66** (2022).

https://doi.org/10.1007/s40194-021-01229-6

[8] L. Na, S. jie Chen, Q. heng Chen, W. Tao, H. Zhao, S. ben Chen, Dynamic welding process monitoring based on microphone array technology. Journal of Manufacturing Processes **64**, 481–492 (2021). https://doi.org/https://doi.org/10.1016/j.jmapro.2020.12.023

[9] N. Lv, Y. Xu, S. Li, X. Yu, S. Chen, Automated control of welding penetration based on audio sensing technology. Journal of Materials Processing Technology **250**, 81–98 (2017). https://doi.org/https://doi.org/10.1016/j.jmatprotec.2017.07.005. URL https://www.sciencedirect.com/science/article/pii/S0924013617302777

[10] A. Sumesh, K. Rameshkumar, K. Mohandas, R.S. Babu, Use of machine learning algorithms for weld quality monitoring using acoustic signature. Procedia Computer Science **50**, 316–322 (2015). https://doi.org/https://doi.org/10.1016/j.procs.2015.04.042. Big Data, Cloud and Computing Challenges

[11] RealityCheck AD https://www.renesas.com/us/en/products/microcontrollers-microprocessors/reality-ai/realitycheck-ad (Last viewed March 27, 2023.)

[12] RealityCheck AD Presentation https://www.youtube.com/watch?v=-6B_XsEN2Q4 (Last viewed March 27, 2023.)

[13] Amazon Monitron https://aws.amazon.com/monitron (Last viewed March 27, 2023.)

[14] Xiris Audio AI https://blog.xiris.com/blog/using-sound-and-imaging-for-detecting-welding-defects (Last viewed March 27, 2023.)

[15] B. Eren, M. Demir, S. Mistikoglu, Recent developments in computer vision and artificial intelligence aided intelligent robotic welding applications. The International Journal of Advanced Manufacturing Technology (2023). https://doi.org/10.1007/s00170-023-11456-4

[16] J. Breitenbach, T. Dauser, H. Illenberger, M. Traub, R. Buettner, in *2021 IEEE International Conference on Big Data (Big Data)* (2021), pp. 2019–2025. https://doi.org/10.1109/BigData52589.2021.9671887

[17] K. Meyer, V. Mahalec, Anomaly detection methods for infrequent failures in resistive steel welding. Journal of Manufacturing Processes **75**, 497–513 (2022). https://doi.org/https://doi.org/10.1016/j.jmapro.2021.12.003

[18] K. Asif, L. Zhang, S. Derrible, E. Indacochea, D. Ozevin, B. Ziebart, Machine learning model to predict welding quality using air-coupled acoustic emission and weld inputs. Journal of Intelligent Manufacturing **33**, 1–15 (2022). https://doi.org/10.1007/s10845-020-01667-x

[19] M. L, K. Senthilkumar, S. Poruran, Performance analysis of weld image classification using modified resnet cnn architecture. Turkish Journal of Computer and Mathematics Education (TURCOMAT) **12**, 2260–2266 (2021)

[20] D. Buongiorno, M. Prunella, S. Grossi, S.M. Hussain, A. Rennola, N. Longo, G. Di Stefano, V. Bevilacqua, A. Brunetti, Inline defective laser weld identification by processing thermal image sequences with machine and deep learning techniques. Applied Sciences **12**(13) (2022). https://doi.org/10.3390/app12136455. URL https://www.mdpi.com/2076-3417/12/13/6455

[21] B. Graney, K. Starry, Rolling element bearing analysis. Materials evaluation **70** (2012)

[22] B. Wang, S. Hu, L. Sun, T. Freiheit, Intelligent welding system technologies: State-of-the-art review and perspectives. Journal of Manufacturing Systems **56** (2020). https://doi.org/10.1016/j.jmsy.2020.06.020

[23] MOTU M4 https://motu.com/en-us/products/m-series/m4 (Last viewed August 3, 2023.)

[24] https://en.wikipedia.org/wiki/Partial_Area_Under_the_ROC_Curve(Last viewed December 7th, 2023.)

[25] Earthworks SR314 https://earthworksaudio.com/vocal-microphones/sr314 (Last viewed August 3, 2023.)

[26] https://xiph.org/flac/features.html(Last viewed November 29, 2023.)

[27] https://en.wikipedia.org/wiki/Partial_Area_Under_the_ROC_Curve(Last viewed November 29, 2023.)

[28] OTC AII V6 http://www.micharc.com/pdfs/otc/OTC%20AII%20Arc%20Robots.pdf (Last viewed August 3, 2023.)

[29] KML Sensors https://www.kmlsensors.com (Last viewed August 3, 2023.)

[30] Slowfast checkpoint https://download.openmmlab.com/mmaction/recognition/slowfast/slowfast_r101_4x16x1_256e_kinetics400_rgb/slowfast_r101_4x16x1_256e_kinetics400_rgb_

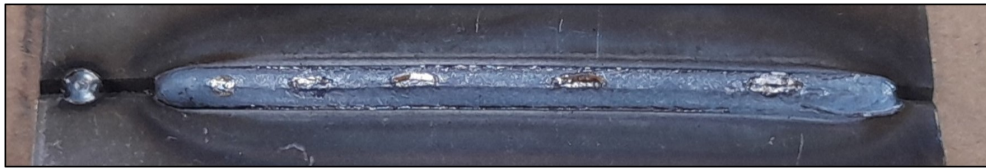20210218-d8b58813.pth (Last viewed November 29, 2023.)

[31] J.A. Lopez, G. Stemmer, P. Lopez Meyer, P. Singh, J. Del Hoyo Ontiveros, H. Cordourier, in *Proceedings of the 6th Detection and Classification of Acoustic Scenes and Events 2021 Workshop (DCASE2021)* (Barcelona, Spain, 2021), pp. 11–15

[32] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, S. Chintala, in *Advances in Neural Information Processing Systems 32*, ed. by H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, R. Garnett (Curran Associates, Inc., 2019), pp. 8024–8035

[33] V.K. Agrawal, S.S. Maurya, Unsupervised detection of anomalous sounds for machine condition monitoring. Tech. rep., DCASE2020 Challenge (2020)

[34] G. Van Rossum, F.L. Drake, *Python 3 Reference Manual* (CreateSpace, Scotts Valley, CA, 2009)

[35] C. Feichtenhofer, H. Fan, J. Malik, K. He, Slowfast networks for video recognition. CoRR **abs/1812.03982** (2018). URL http://arxiv.org/abs/1812.03982

[36] M. Contributors. Openmmlab's next generation video understanding toolbox and benchmark. https://github.com/open-mmlab/mmaction2 (2020)

[37] D.P. Kingma, J. Ba. Adam: A method for stochastic optimization (2017)

[38] L.N. Smith, N. Topin. Super-convergence: Very fast training of neural networks using large learning rates (2018)

[39] T. Fawcett, An introduction to roc analysis. Pattern Recogn. Lett. **27**(8), 861–874 (2006)
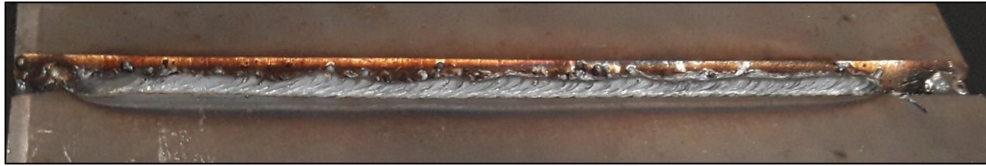
(a) Good, Normal-condition, Sample.


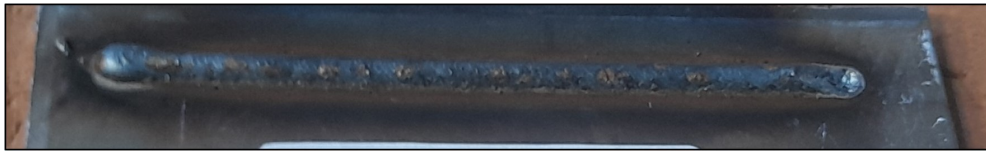(b) Burnthrough Sample.


(c) Excessive Penetration Sample.


(d) Lack Of Fusion Sample.


(e) Porosity Sample.


(f) Spatter Sample.


(g) Warping Sample.

Fig. B1: Welding Samples