

# WHITED - A Worldwide Household and Industry Transient Energy Data Set

Matthias Kahl, Anwar Ul Haq, Thomas Kriechbaumer and Hans-Arno Jacobsen  
Technische Universität München  
Email: matthias.kahl@in.tum.de

**Abstract**—In this paper, we introduce a data set of appliance start-up measurements from several locations. The appliances were recorded with a low-cost custom sound card meter. The recording was mainly done in households and small industry settings in different regions around the world. Thus, it may be possible to extract region-specific grid characteristics from the voltage waveforms in the data. To cover all corresponding transients, we recorded the first 5 seconds of the appliance start-ups for 110 different appliances to date, amounting to 47 different appliance types. The aim of this data set is to provide a broad spectrum of different appliance types in regions around the world.

## I. INTRODUCTION

A trend towards the incorporation of green energy technologies into the existing grid is on the rise. Much of the renewable energy is generated and utilized at the consumer end. For efficient integration of highly weather dependent distributed resources, consumers need to play an effective role. Also, they need to be incentivized by providing them with information of their real-time energy consumption, preferably at the level of individual appliances. So the importance of handling energy consumption is growing and has even lead to new fields of research like non-intrusive load monitoring (NILM) [1]. NILM relies on advances in computer science by realizing abstractions to help retrieve knowledge from energy consumption data. Applications are found in energy measurement scenarios including household energy demand, appliance transients, and reduction of electro-magnetic interferences (EMI) [2]–[10]. Once such abstractions are recorded, they constitute the data sets that provide useful information about the consumer behavior by breaking down energy consumption for each appliance.

Although NILM has been around for almost three decades, it was not until recently that this field started to flourish. Apart from recent advances in machine learning techniques, one of the major contribution was the release of an open energy data set REDD [3]. Although many open data sets have been released since then, we believe that there is still a scarcity of publicly available energy data sets in the high frequency domain. As stated by K. C. Armel et al. in [11], the high frequency sampling rate in electricity signal measurements enables us to accurately distinguish between more appliances. Unfortunately, most of the time, high frequency data acquisition hardware for transient analysis is expensive. In the UK-DALE data set, the authors have demonstrated the use of an on-board sound card to record the electricity signals

TABLE I  
COMPARISON OF DATA SETS WITH HIGH FREQUENCY APPLIANCE TRACES

Data Set	Bit	Fs	Appliance		Purpose
			Classes	Variety	
REDD [5]	24	15 kHz	~ 20	10	house demand
BLUED [2]	16	12 kHz	~ 30	~ 1	house demand
UK DALE [4]	20	16 kHz	~ 40	~ 1 – 3	house demand
PLAID [8]	16	30 kHz	~ 12	~ 20	transients
HFED [7]	16	5 MHz	15	1	spectral traces
WHITED	16	44 kHz	46	1 – 9	transients

at high frequency [4]. Sound cards have been around for years with high resolution analog to digital converters. Many lossless compression algorithms are available for audio data to reduce the data size. For our measurements, we used the stereo line input of an external USB sound card, since most modern laptops lack an dedicated stereo line-in port.

We have collected a general purpose and freely accessible data set, using an affordable and portable sound card measurement system. This paper presents the sound card measurement system based on simple assembly instructions, and it demonstrates several characteristics of the collected data set. In future work, we plan to extend the data set with additional measurements by drawing on crowd sourcing mechanisms. One of our main insights is that even low cost hardware perform pretty well for NILM purposes.

## II. RELATED WORK

Several public data sets covering appliance-level energy consumption already exist. The purpose of these data sets is to measure demand in private households through a non-intrusive single point measurement in either low or high frequency. Through constant observation of household energy demand, these data sets provide comprehensive longtime measurements to cover user behavior in the corresponding residence. These data sets are a good source for power disaggregation tasks as they indirectly provide transient start-up features at an appliance level. Real-world scenario data sets include REDD [5], UK DALE [4] and BLUED [2] among others.

When looking in more detail at appliance transients, it can be cumbersome to extract them from these single measurements. Since the ground truth is mostly based on 1s to 6s data without explicit voltage or current waveforms, it might be possible that two start-ups fall in the same time window, thus,

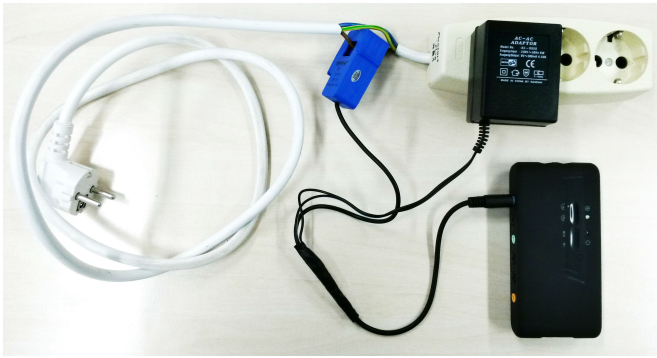


Fig. 1. Measurement equipment prototype

violating the assumption of the switch continuity principle (SCP) [12]. Therefore, it is helpful to take a closer look at transient-focused data sets such as PLAID [8] and HFED [7]. PLAID examines start-up transients at 30 kHz whereas in HFED short transient spectral traces of up to 5 MHz were observed but require high effort in terms of hardware and experimental setup to reproduce.

Table I gives a comparison between the above mentioned high frequency data sets in terms of resolution, purpose, amount of appliance types (classes) and quantity of appliances for each class (Variety). The information about the appliance types and quantity are inferred from the available data. We believe that a high intra-class variety leads to a more reliable result in terms of appliance classification.

With WHITED – a Worldwide Household and Industry Transient Energy Data set – we want to contribute to existing energy data sets in terms of higher sampling frequency and higher amount of appliance types and variety. In addition, we provide a region classification for each measurement to potentially enable the investigation of region specific research questions.

### III. ARCHITECTURE

In this section, we describe hard and software components of our measurement equipment which is based on a sound card as inexpensive analog to a digital converter. The idea of a sound card-based measurement system is not new and was already used in [4] and [13]. Sound cards have a very good price vs. performance ratio when using them as an analog to digital converter. Our measurement prototype is based on a modified 3-port extension cord, a current clamp, an AC-AC transformer, a voltage divider, and an external USB sound card with a *Cmedia CM6206* chipset.

#### A. Hardware Design

For measuring the current, we use a YHDC current clamp with built-in burden resistor. This current clamp produces a 1 V signal at 30 A primary current. For the voltage measurements, we need to transform the grid voltage from 230 V to 11 V with the AC-AC transformer. To have a corresponding voltage signal that lies in the line-in range of the sound card, we reduce it with a voltage divider to 0.47 V. The voltage

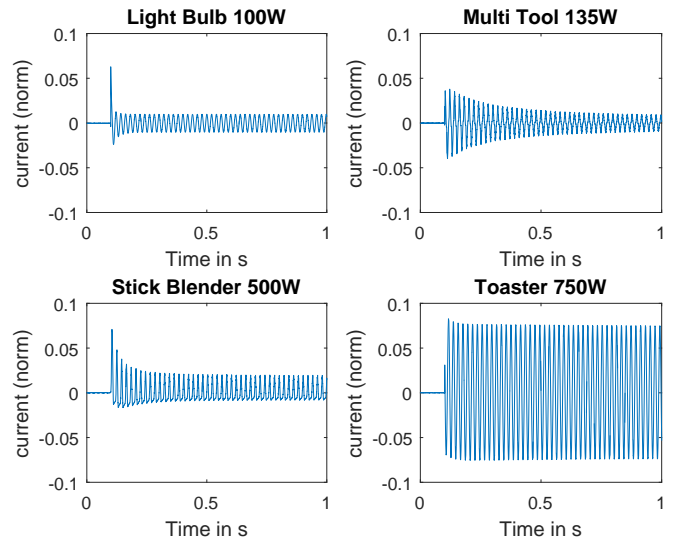


Fig. 2. Start-up of four different appliances. The different in-rush current characteristics are clearly visible.

divider is located in the black isolation part that merges the current and voltage signal cables into one cable that goes into the sound card. See Fig. 1 for the complete configuration.

#### B. Measurements

The signals were recorded in 44.1 kHz temporal and 16 bit amplitude resolution. To be able to take multiple measurements in different places, it was necessary to build 3 identical measurement kits. Therefore, we also have to deal with three slightly different sets of calibration factors. The calibration itself is done with an *VOLTCRAFT VC-330* multimeter. Since the multimeter provides current measurements with a current clamp, it was possible to measure both signals - voltage and current - and define a voltage and current calibration factor for each measurement kit. Some sample measurements can be seen in Fig. 2 for 4 different appliances.

To cover the start-up transients of the appliances, it is necessary to determine them on demand. This is implemented with a *Matlab* routine that uses the internal DSP package to monitor the line-in signal of the sound card. The start-up is defined based on the current signal energy crossing a threshold. If the current signal energy leads to a start-up, the routine starts recording and adds 100 ms of the signal beforehand as pre-start-up window. This window allows difference-based algorithms to work effectively. That means that not the absolute power consumption on the start-up but the difference between the power of the pre-start-up window and the start-up power can be observed. This approach introduces more flexibility for developing algorithms that allow the recognition of concurrently running appliances with different start-up transients.

We decided to measure 10 start-ups for each appliance. These start-ups were triggered manually by the user. Appliances that have no switch (e.g., an iron) were just plugged and unplugged 10 times as it would be the case under real

usage. The appliances are measured for 5 seconds which is the duration of each start-up we recorded.

### C. Data Set

To this end, our data set comprises 1100 different records for 110 different appliances which can be grouped into 47 different types (classes) in 6 different regions. For most appliances, we took a photo of its electrical specification label. These images are located in the sub-folder `images` and `type-labels`. Table II gives an overview of the measured appliances. The signal containing files are saved as `flac` files – a common lossless audio file format. The file names contain meta information and are of following format:

```
[Class]_[Name]_[Region]_[#Kit]_[TimeStamp].flac
GuitarAmp_Marshall18240_R3_MK2_20151115133402.flac
```

The data set is freely available on the following web page: <http://bit.ly/WHITED-Set>. For demand, load and appliance information retrieval, the most important signal is the current. To give the voltage signal a higher significance, we decided to measure the voltage in several regions that follow the European grid standards. To this end, the data set contains 4 regions in Germany, 1 in Austria, and 2 in Indonesia.

Since grid characteristics are mainly affected by utilities and the consumption characteristics of the surrounding area, a future research direction is to look for possibilities to determine the region from the voltage signal. This experiment is a similar classification task to the appliance recognition we have already implemented.

## IV. EVALUATION

To ensure the quality of the data set, we applied several signal quality checks and conducted two classification experiments.

### A. Data Quality

Since sound cards do not provide a high level of linearity in frequency response as compared with professional ADCs, we verified that there is no significant impact on the measurements taken.

The sound card manufacturer provides some information regarding the line-in linearity which can be seen in Fig. 4. It is visible that the strongest damping of around 0.25 dB has its maximum at 3320 Hz. The steepest flank has a bandwidth of around 3300 Hz and lies between 3320 Hz and 6622 Hz which is acceptable for most considered purposes.

To obtain an approximation of the noise level during recording, the energy of a 10 second empty signal is being compared to the energy of a maximum amplitude sine-wave signal. With this calculation, we estimate an effective SNR (signal to noise ratio). We measured an average noise RMS of 4.8 mA where 30 A corresponds to the RMS maximum.

$$SNR = 20 \cdot \log_{10} \frac{RMS_{max}}{RMS_{noise}}$$

TABLE II  
APPLIANCE TYPES (CLASSES) THAT WERE MEASURED

AC	1	Air Pump	1	Bench Grinder	1
CFL	2	Charger	7	Coffee Machine	1
Deep Fryer	1	Desktop PC	1	Desoldering tool	1
Drilling Machine	2	Fan	6	Fan Heater	1
Flat Iron	2	Game Console	4	Guitar Amp	1
Hair Dryer	6	Halogen Fluter	1	Heater	1
HiFi Rack	1	Iron	3	Jigsaw	1
JuiceMaker	1	Kettle	6	Laptop	1
Laserprinter	1	LED Light	9	Light bulb	6
Massage tool	3	Microwave	2	Mixer	4
Monitor	2	Mosquito Repellent	1	Multitool	1
Powersupply	4	Projector	1	Sewing Machine	1
Shoe warmer	2	Shredder	2	Soldering Iron	2
Toaster	4	Treadmill	1	TV	1
Vacuum Cleaner	4	Washing Machine	1	Water Heater	4
Water Pump	1				

$$SNR = 20 \cdot \log_{10} \frac{30 \text{ A}}{0.0048 \text{ A}} = 75.91 \text{ dB}$$

The effective SNR of this measurement system is 75.91 dB. The maximum measurable peak to peak current  $I_{p-p}$  is  $30.0 \text{ A}_{RMS} \cdot 2\sqrt{2} = 84.4 \text{ A}$ . Therefore, we calculate an effective current resolution with a step size of 13.5 mA.

$$I_{step} = \frac{I_{p-p}}{I_{maxRMS}} \cdot I_{noiseRMS}$$

$$I_{step} = \frac{84.4 \text{ A}}{30.0 \text{ A}} \cdot 0.0048 \text{ A} = 0.0135 \text{ A}$$

This current step size enables us to calculate the effective power step size  $P_{step}$  corresponding to 230 V of grid voltage.

$$P_{step} = 230 \text{ V} \cdot 0.0135 \text{ A} = 3.1 \text{ W}$$

The resolution and noise of the sound card allows a voltage step of 0.313 V, a current step of 0.0135 A which results in a measurable power step of around 3.1 W based on 230 V. To achieve reliable results only appliances with a consumption of at least 20 W are considered in our data set. This covers most household and small industry appliances.

Fig. 3a and 3b show a spectrogram of a mixer and a multi-tool based on the first 5 seconds after the start-up. Both appliances have a fast spinning motor and look similar in the time domain. However, there are significant differences in the spectral domain that can be transformed into distinguishable features for appliance classification purposes.

### B. Experimental Results

Our appliance recognition experiment is based on a classification task to distinguish appliances on its characteristics in the current signal. The classifier has to distinguish between all 47 appliance types. The classification experiment is implemented in Matlab. All `flac` files are imported and the containing

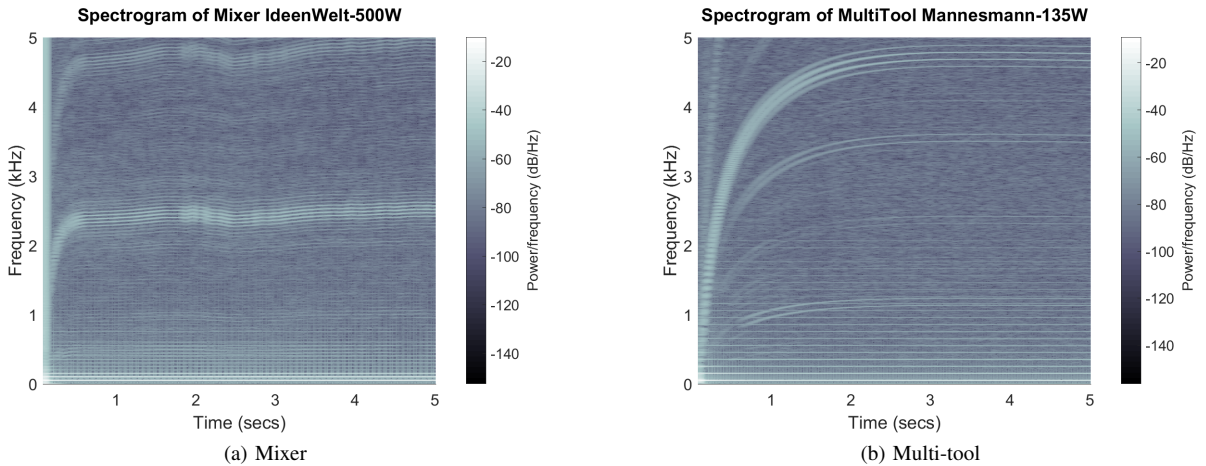


Fig. 3. Comparison of 2 appliances that use motors with relatively high rotations per minute. The different spectral characteristics are clearly visible. The multi-tool has stronger uneven harmonics while the harmonics are more equal in the case of the mixer.

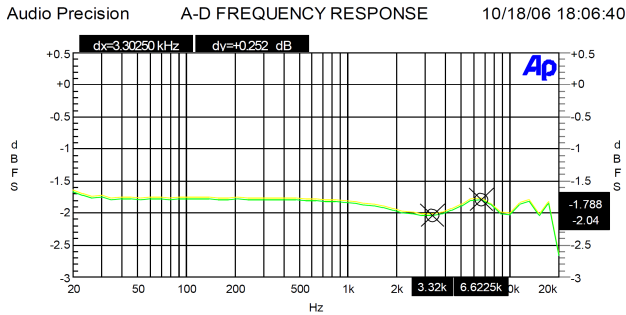


Fig. 4. The line-in frequency response from the CM6206 specification [14].

signal is scaled with the corresponding calibration factors to determine actual values. After this preprocessing step, a region of interest (ROI) needs to be extracted. Here, we decided to cut the signal right on the start-up until 500 ms after the start-up. These 500 ms samples are given to the feature extraction stage which is an implementation of 13 different characteristics including harmonics, phase shift and total harmonic distortion (THD).

The best results we achieved for the appliance classification were based on a feature set that consisted of a period-based power trend with 25 dimensions, the THD and crest factor of the current spectrum with each 1 dimension in its size. With these three features in 27 dimensions, we achieve an average classification accuracy across all appliances of around 95 % with a 10-fold cross-validation and a support vector machine (SVM) classifier. This confirms the observation that power difference and harmonics contain sufficient information to distinguish among basic electrical appliances [15].

For the region classification experiment, we use the same environment but employ the voltage instead of current for the feature extraction. The labels are not the appliances but the region where the measurements were taken. We apply the voltage, grid frequency and a few spectral- and waveform-

based features. We obtain an almost perfect classification accuracy of 99.13 % with an SVM classifier. Here, we must consider that the feature extraction is based on characteristics that vary over time and are not independently representative for the corresponding region.

## V. CONCLUSIONS

In this work, we publish a data set comprised of a broad range of household and small industry appliance start-up transients. As discussed, we believe that there is still a need for such kind of measurements. The purpose of this paper is to show that even a low-budget, custom measurement system allows one to retrieve significantly discriminating features from appliance start-up transients to enable appliance classification needs.

We aim at continuing to expand our data set through involving the community and hopefully more individuals world-wide join to contribute measurements based on the measurement system specification. Location recognition based on the voltage signal needs further observations of grid characteristics in each region to be able to distinguish between regions. The reason is that grid characteristics like frequency and voltage constantly changing. Short duration measurements do not provide sufficient information about the stability of such grid characteristics.

## REFERENCES

- [1] G. W. Hart, "Nonintrusive appliance load monitoring," *Proceedings of the IEEE*, vol. 80, no. 12, pp. 1870–1891, 1992.
- [2] K. Anderson, A. Ocneanu, D. Benitez, D. Carlson, A. Rowe, and M. Berges, "Blued: A fully labeled public dataset for event-based non-intrusive load monitoring research," in *Proceedings of the 2<sup>nd</sup> KDD workshop on data mining applications in sustainability (SustKDD)*, 2012, pp. 1–5.
- [3] C. Beckel, W. Kleiminger, R. Cicchetti, T. Staake, and S. Santini, "The eco data set and the performance of non-intrusive load monitoring algorithms," in *Proceedings of the 1<sup>st</sup> ACM Conference on Embedded Systems for Energy-Efficient Buildings*, 2014, pp. 80–89.

- [4] J. Kelly and W. Knottenbelt, "The uk-dale dataset, domestic appliance-level electricity demand and whole-house demand from five uk homes," 2015. [Online]. Available: <http://www.doc.ic.ac.uk/~dk3810/data/>
- [5] J. Z. Kolter and M. J. Johnson, "Redd: A public data set for energy disaggregation research," in *Workshop on Data Mining Applications in Sustainability (SIGKDD)*, San Diego, CA, vol. 25. Citeseer, 2011, pp. 59–62.
- [6] A. Monacchi, D. Egarter, W. Elmenreich, S. D'Alessandro, and A. M. Tonello, "Greend: An energy consumption dataset of households in Italy and Austria," *CoRR*, vol. abs/1405.3100, 2014.
- [7] M. Gulati, S. Sundar Ram, and A. Singh, "An in depth study into using EMI signatures for appliance identification," in *Proceedings of the First ACM International Conference on Embedded Systems For Energy-Efficient Buildings*. ACM, 2014.
- [8] J. Gao, S. Giri, E. C. Kara, and M. Bergés, "Plaid: a public dataset of high-resolution electrical appliance measurements for load identification research: demo abstract," in *Proceedings of the 1<sup>st</sup> ACM Conference on Embedded Systems for Energy-Efficient Buildings*. ACM, 2014, pp. 198–199.
- [9] A. Veit, C. Goebel, R. Tidke, C. Doblender, and H.-A. Jacobsen, "Household electricity demand forecasting: Benchmarking state-of-the-art methods," in *5th International Conference on Future Energy Systems*, ser. e-Energy '14. New York, NY, USA: ACM, 2014, pp. 233–234.
- [10] H. Ziekow, C. Doblender, C. Goebel, and H.-A. Jacobsen, "Forecasting household electricity demand with complex event processing: Insights from a prototypical solution," in *13th ACM/FIP/USENIX International Middleware Conference*, ser. Middleware Industry'13. New York, NY, USA: ACM, 2013, pp. 2:1–2:6.
- [11] K. C. Armel, A. Gupta, G. Shrimali, and A. Albert, "Is disaggregation the holy grail of energy efficiency? the case of electricity," *Energy Policy*, vol. 52, pp. 213–234, 2013.
- [12] S. Makonin, "Investigating the switch continuity principle assumed in non-intrusive load monitoring (nilm)," 2016.
- [13] F. Englert, T. Schmitt, S. Köbler, A. Reinhardt, and R. Steinmetz, "How to auto-configure your smart home?: high-resolution power measurements to the rescue," in *Proceedings of the fourth international conference on Future energy systems*. ACM, 2013, pp. 215–224.
- [14] *CM6206 High Integrated USB Audio I/O Controller*, Rev. 2.1 ed., C-Media Electronics Inc. [Online]. Available: <http://www.bramcam.nl/NA/8663-XS/CM6206.pdf>
- [15] K. N. Trung, O. Zammit, E. Dekneuveel, B. Nicolle, C. N. Van, and G. Jacquemod, "An innovative non-intrusive load monitoring system for commercial and industrial application," in *Advanced Technologies for Communications (ATC), 2012 International Conference on*. IEEE, 2012, pp. 23–27.