# FRIEDRICH-SCHILLER-UNIVERSITÄT JENA

# Data Science Day Jena 2023

May 10th, 2023 – Rosensäle of the Friedrich Schiller University Jena

# Book of Abstracts

Julien Klaus, Olia Blacher (Ed.)
{julien.klaus, olia.blacher}@uni-jena.de

# Contents

# Potential and limits of minimal informative priors for hidden Markov Models to improve parameter inference

**Author:** Jan Münch[1]

**Co-author:** Klaus Benndorf [1]

[1] *Physiology 2, Friedrich Schiller University Jena*

**Corresponding Author:** jan.muench@med.uni-jena.de

Inferring the complex functional dynamics of ion channels from ensemble currents is a daunting task. We previously addressed this problem by applying a parallelized Bayesian filter to specify kinetic schemes for macroscopic current and fluorescence data leeding to a more accurate likelihood than previous gold-standard algorithms (Muench 2022 eLife 11:e62714).

Using Bayesian statistics requires to define a prior distribution. When little information about the parameter is known, especially when the information content in the data is poor, the prior is crucial to make the posterior as sensitive as possible to the data. For ion-channel HMMs, a minimal informative prior may consist of a log uniform prior for the inverse dwell times of a state and a Dirichlet prior for the probability of each transition out of the state. Applying this prior reduces the number of ion channels required for a reasonable inference by one order of magnitude compared to the standard uniform prior, which is often considered mistakenly as uninformative.

Ion currents from patch-clamp experiments observe only partially the dynamics of interest in the chemical network. We show by simulated patch-clamp data that this partial observability causes the likelihood to become flat in many directions in the parameter space, causing a degree of practical parameter non-identifiability for most non-trivial hidden Markov models. The log uniform distribution of the minimal informative prior desensitizes the posterior to the non-identifiability problem of the likelihood for some part of the parameter space. Nevertheless, the posterior will always be dominated by the structure of the prior in the rest of the parameter space. Thus, all posteriors of HMM models will be improper if equipped with minimal informartive improper prior distributions. Here, we discuss how to recognize and treat this practical parameter non-identifiabilty problem, inherent in most HHMs and any statistical framework used. Including other physically constraints further increases the inference quality and decrease the parameter unidentifiabilty problem. We conclude that for a given data quality the quality of model inference can be significantly improved by selecting a minimal informative prior.

# Predicting anharmonicity constants using machine learning

**Author:** Jamoliddin Khanifaev[1]

**Co-author:** Eva von Domaros

[1] *Friedrich Schiller University Jena*

**Corresponding Author:** jamoliddin.khanifaev@uni-jena.de

Considering anharmonic effects is essential to accurately describe vibrational, spectroscopic and thermodynamic properties of molecular systems. Recently, it has been shown that introducing anharmonic formalism into the quantum cluster equilibrium (QCE) program improves the results of the calculations [1]. It is possible to calculate anharmonic frequencies [2]. However, such calculations are tedious and computationally expensive, therefore it is of prime interest to apply machine learning techniques to overcome this task. In this work we study a variety of molecular clusters of different sizes consisting of HX, CH3X, C2H5X (X = F, Cl, Br) monomers. Our dataset features consist of normal mode coordinates as well as harmonic frequencies, anharmonic frequencies and intensities of the fundamental and first overtone modes. Symmetry and structural descriptors such as internal coordinates are also taken into account. Anharmonicity constants are extracted from the vibrational energy levels of the Morse oscillator. Later on, various machine learning algorithms are applied for the classification and regression purposes. Our results show that while stretching vibration modes have positive anharmonicity constants, the low frequency bending and torsional modes often exhibit negative values. With this we were able to classify the anharmonicity constants according to the type of vibration such as stretching, bending, or internal rotations and translations. Regression algorithms were able to provide an estimate of the anharmonic frequencies.

### References

[1] J. Chem. Phys. 146, 124114 (2017)

[2] J. Chem. Phys. 105, 10332 (1996)

# The Archive Query Log: Mining Millions of Search Result Pages of Hundreds of Search Engines from 25 Years of Web Archives

**Authors:** Jan Heinrich Reimer[1]; Sebastian Schmidt[2]; Maik Fröbe[1]; Lukas Gienapp[2]; Harrisen Scells[2]; Benno Stein[3]; Matthias Hagen[1]; Martin Potthast[2,4]

[1] *Friedrich Schiller University Jena*

[2] *Leipzig University*

[3] *Bauhaus-Universität Weimar*

[4] *ScaDS.AI*

**Corresponding Author:** heinrich.reimer@uni-jena.de

The Archive Query Log (AQL) is a previously unused, comprehensive query log collected at the Internet Archive over the last 25 years. Its first version includes 356 million queries, 166 million search result pages, and 1.7 billion search results across 550 search providers. Although many query logs have been studied in the literature, the search providers that own them generally do not publish their logs to protect user privacy and vital business data. Of the few query logs publicly available, none combines size, scope, and diversity. The AQL is the first to do so, enabling research on new retrieval models and (diachronic) search engine analyses. Provided in a privacy-preserving manner, it promotes open research as well as more transparency and accountability in the search industry.

# Simulation study for artifacts markers in microscopic images

**Authors:** Elena Corbetta[1]; Thomas Bocklitz[2]

[1] *Institute of Physical Chemistry (IPC) and Abbe Center of Photonics (ACP), Friedrich Schiller University Jena, Member of the Leibniz Centre for Photonics in Infection Research (LPI)*

[2] *Leibniz Institute of Photonic Technology, Member of Leibniz Health Technologies, Member of the Leibniz Centre for Photonics in Infection Research (LPI)*

**Corresponding Author:** elena.corbetta@uni-jena.de

Optical microscopy is a powerful and minimally invasive tool for the investigation of biological processes. In this context, processing of images is of utmost importance to improve image quality and sample understanding.[1] However, there is not a standard quantitative approach to evaluate image quality, especially in presence of artifacts. The computation of metrics can provide ambiguous results, with poor agreement of the metrics with human visual perception. In addition, the ground truth is often needed for comparison.[2] To address these issues, we performed a systematic study to identify markers and metrics for the characterization and evaluation of microscopic images. We developed simple models for simulation of biological structures and the most common microscopic artifacts; these include blurring, mixed Poisson-Gaussian noise, and uneven illumination. The models can be applied by tuning independent parameters to modulate the sample structure or the specific effect of the artifact. The metrics for image evaluation were selected after extensive literature research, taking as reference previous studies on microscopic measurements. [2][3] We obtained a collection of images with a variety of simulated experimental conditions and specific trends of the metrics were identified for each artifact, developing an overview of reference markers for different degradations. Finally, image markers were validated on real experimental datasets. These results help the understanding of experimental acquisitions and should be considered when evaluating the effect of different processing workflows on the same input image.

Figure 1: Schematics of the workflow for evaluation of image markers. The left section shows a comparison between simulations (column (a)) and real images (column (b)) of three of the selected biological structures. The central section displays different artifacts applied to the (c) simulated and (d) experimental images of cell nuclei; from the top: blurring, combination of dark, shot and readout noise modeled as mixed Poisson-Gaussian noise, and uneven illumination. Section (e) shows a schematic representation of the evaluation of image markers, with a qualitative example plot that can be obtained when comparing the same metrics on different artifacts.

## References

[1] J. Roels et al., Adv Anat Embryol Cell Biol. 219 (2016), 41-67

[2] S. Koho et al., Sci Rep 6 (2016), 28962

[3] R.P.J. Nieuwenhuizen et al., Nat Methods 10 (2013), 557-562

# Causal Discovery using Model Invariance through Knockoff Interventions

**Author:** Wasim Ahmad[1]

**Co-authors:** Joachim Denzler[1]; Maha Shadaydeh[1]

[1] *Computer Vision Group, Friedrich Schiller University Jena*

**Corresponding Author:** wasim.ahmad@uni-jena.de

Cause-effect analysis is crucial to understand the underlying mechanism of a system. We propose to exploit model invariance through interventions on the predictors to infer causality in nonlinear multivariate systems of time series. We model nonlinear interactions in time series using DeepAR and then expose the model to different environments using Knockoffs-based interventions to test model invariance. Knockoff samples are pairwise exchangeable, in-distribution and statistically null variables generated without knowing the response. We test model invariance where we show that the distribution of the response residual does not change significantly upon interventions on non-causal predictors. We evaluate our method on real and synthetically generated time series. Overall our method outperforms other widely used causality methods, i.e, VAR Granger causality, VARLiNGAM and PCMCI+. The code and data can be found at: https://github.com/wasimahmadpk/deepCausality

# Why Capsule Neural Networks Do Not Scale

**Authors:** Matthias Mitterreiter[1]; Marcel Koch[1]; Joachim Giesen[1]; Sören Laue[2]

[1] *Friedrich Schiller University Jena*

[2] *University of Hamburg*

**Corresponding Author:** matthias.mitterreiter@uni-jena.de

Capsule neural networks replace simple, scalar-valued neurons with vector-valued capsules. They are motivated by the pattern recognition system in the human brain, where complex objects are decomposed into a hierarchy of simpler object parts. Such a hierarchy is referred to as a parse-tree. Conceptually, capsule neural networks have been defined to realize such parse-trees. The capsule neural network (CapsNet), by Sabour, Frosst, and Hinton, is the first actual implementation of the conceptual idea of capsule neural networks. CapsNets achieved state-of-the-art performance on simple image recognition tasks with fewer parameters and greater robustness to affine transformations than comparable approaches. This sparked extensive follow-up research. However, despite major efforts, no work was able to scale the CapsNet architecture to more reasonable-sized datasets. Here, we provide a reason for this failure and argue that it is most likely not possible to scale CapsNets beyond toy examples. In particular, we show that the concept of a parse-tree, the main idea behind capsule neuronal networks, is not present in CapsNets. We also show theoretically and experimentally that CapsNets suffer from a vanishing gradient problem that results in the starvation of many capsules during training.

# Optimization for Machine Learning

**Authors:** Matthias Mitterreiter[1]; Sören Laue[2]; Joachim Giesen[1]

[1] *Friedrich Schiller University Jena*

[2] *University of Hamburg*

**Corresponding Author:** matthias.mitterreiter@uni-jena.de

Optimization is an integral part of most machine learning systems and most numerical optimization schemes rely on the computation of derivatives. Therefore, frameworks for computing derivatives are an active area of machine learning research. Surprisingly, as of yet, no existing framework is capable of computing higher order matrix and tensor derivatives directly. Here, we close this fundamental gap and present an algorithmic framework for computing matrix and tensor derivatives that extends seamlessly to higher order derivatives. The framework can be used for symbolic as well as for forward and reverse mode automatic differentiation. Experiments show a speedup of up to two orders of magnitude over state-of-the-art frameworks when evaluating higher order derivatives on CPUs and a speedup of about three orders of magnitude on GPUs.

But, although optimization is the longstanding, algorithmic backbone of machine learning new models still require the time-consuming implementation of new solvers. As a result, there are thousands of implementations of optimization algorithms for machine learning problems. A natural question is, if it is always necessary to implement a new solver, or is there one algorithm that is sufficient for most models. Common belief suggests that such a one-algorithm-fits-all approach cannot work, because this algorithm cannot exploit model specific structure. At least, a generic algorithm cannot be efficient and robust on a wide variety of problems. Here, we challenge this common belief. We have designed and implemented the optimization framework GENO (GENeric Optimization) that combines a modeling language with a generic solver. GENO takes the declaration of an optimization problem and generates a solver for the specified problem class. The framework is flexible enough to encompass most of the classical machine learning problems. We show on a wide variety of classical but also some recently suggested problems that the automatically generated solvers are (1) as efficient as well engineered, specialized solvers, (2) more efficient by a decent margin than recent state-of-the-art solvers, and (3) orders of magnitude more efficient than classical modeling language plus solver approaches.

# Understanding the effects of plant diversity on soil temperature stability

**Authors:** Gideon Henrik Stein[1,2]; Yuanyuan Huang[2]

**Co-authors:** Maha Shadaydeh[1]; Anne Ebeling[1]; Nico Eisenhauer; Joachim Denzler[1]

[1] *Friedrich Schiller University Jena*

[2] *iDiv, University Leipzig*

**Corresponding Author:** gideon.stein@uni-jena.de

To understand the mechanism that drives the relationship between plant diversity and soil temperature stability, we compare the results of Structural Equation Modeling (SEM) based on domain-expert hypotheses with those of two Causal Inference methods which are entirely agnostic towards possible relationship (The PC algorithm [1] and the CCDr method [2]). While the results are generally consistent, we also observe disparities between methods.

## References

[1] Colombo, Diego, and Marloes H. Maathuis. "Order-independent constraint-based causal structure learning." J. Mach. Learn. Res. 15.1 (2014): 3741-3782.

[2] Aragam, Bryon, and Qing Zhou. "Concave penalized estimation of sparse Gaussian Bayesian networks." The Journal of Machine Learning Research 16.1 (2015): 2273-2328.

# Pre-processing Raman data via deep learning method

**Author:** Azadeh Mokari[1,2]

**Co-authors:** Simone Eiserloh[1]; Ute Neugebauer[1,2]; Thomas Bocklitz[1,2,3]

[1] *Leibniz Institute of Photonic Technology, Member of Leibniz Health Technologies, Member of the Leibniz Centre for Photonics in Infection Research (LPI), Albert-Einstein-Strasse 9, 07745 Jena, Germany.*

[2] *Institute of Physical Chemistry (IPC) and Abbe Center of Photonics (ACP), Friedrich Schiller University Jena, Member of the Leibniz Centre for Photonics in Infection Research (LPI), Helmholtzweg 4, 07743 Jena, Germany.*

[3] *Institute of Computer Science, Faculty of Mathematics, Physics & Computer Science, University Bayreuth Universitaetsstraße 30, 95447 Bayreuth, Germany.*

**Corresponding Author:** azadeh.mokari@uni-jena.de

Raman spectroscopy is a type of analytical technique that uses the interaction of light with a sample to provide information about its atomic and molecular properties. However, Raman spectra are frequently overshadowed by inconsistencies in baselines and various sources of noise. These defects and contributions to the Raman data must be rectified before identifying or categorizing the samples. Accordingly, Raman data is processed using AI-based algorithms. To that end, we suggested the use of a deep learning approach as a pre-processing tool for Raman data. As a result, we tested two networks: the convolutional denoising autoencoder (CDAE) [1] and U-Net [2]. CDAE and U-Net networks were implemented to test two different pre-processing cases: denoising and denoising with baseline removal. In both cases, the superiority of the methods was evaluated using real and artificial Raman data. In the first case, we aimed to reconstruct high-quality (HQ) Raman spectra that included a background. Therefore, the networks were trained to map between noisy Raman data measured with different integration times, for example, 0.5 s as an input and HQ Raman data with 1 s as an output. As shown in Figure 1, the U-Net/CDAE network tries to estimate the HQ data in experiment data or predict the HQ artificial Raman data. Afterward, in the testing phase, the trained networks are used to predict the HQ data. In the second case, we aimed to reconstruct high-quality spectra with baseline removal. In other words, the aim of this case is to remove noise and background from the data at the same time. Therefore, the same noisy Raman data was used as an input, and the output was acquired by applying classical pre-processing methods (SG+SNIP on the HQ Raman data). Regarding the evaluation part in Figure 1, U-Net has the capability to remove the noise and baseline simultaneously, while the CDAE is only able to remove the noise. In conclusion, the suggested technique outperforms traditional methods in terms of time and error.

## Acknowledgements

## References

[1] Fan, X.g., et al., Signal-to-noise ratio enhancement for Raman spectra based on optimized Raman spectrometer and convolutional denoising autoencoder. Journal of Raman Spectroscopy, 2021. 52(4): p. 890-900.

[2] Guo, S., et al., Deep learning for 'artefact' removal in infrared spectroscopy. Analyst, 2020. 145(15): p. 5213-5220.
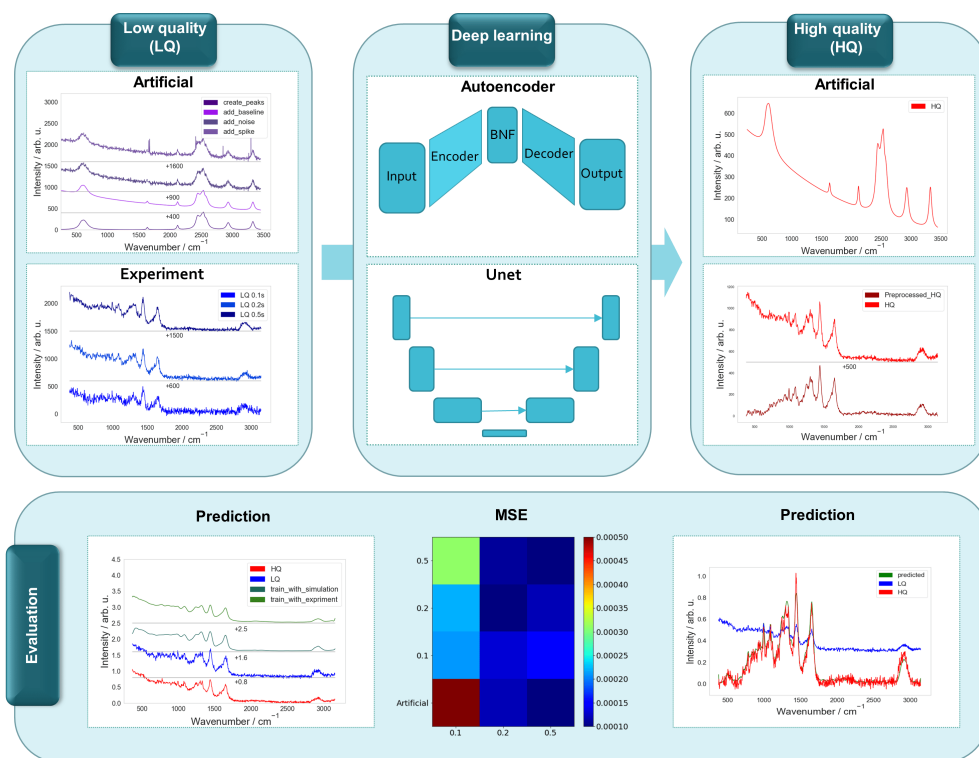
Figure 1: In the reconstruction of the HQ spectra, the training phase includes noisy data with different integration times, for example, 0.5s as an input and Raman data with 1s in experiment data as an output. Therefore, the network learns to map the noisy signal to the output signal acquired in a longer period. In the testing phase, the noisy data is reconstructed according to the networks and the model. In the reconstruction of HQ spectra with baseline removal, the HQ data is reconstructed by classical pre-processing methods (Savitzky-Golay and baseline correction) in order to reach the predicted spectra without noise and baseline simultaneously. In the evaluation part, the result of prediction by two types of networks was plotted. In addition, MSE shows the error in the different possibilities of training and testing with different integration times.

# Fine-grained Recognition and Continuous Learning for Biodiversity Monitoring

**Authors:** Daphne Frederike Auer[1]; Dimitri Korsch; Joachim Denzler; Julia Böhlke[1]; Matthias Frank Körschens[1]; Paul Bodesheim[1]

[1] *Friedrich Schiller University Jena*

**Corresponding Author:** dimitri.korsch@uni-jena.de

Biodiversity monitoring is crucial to understand and preserve our nature. Hence, the computer vision group supports ecologists with AI-based tools for the automated evaluation of images. First, we develop methods for fine-grained species recognition, which is required since the visual appearance of related species complicates traditional recognition methods. We also present methods to enrich the training data with freely available but noisy annotated images from the Internet. Furthermore, our continuous learning methods allow us to quickly adapt trained models to new data and changing requirements with low manual annotation effort. Finally, besides animals, we also develop methods for plant-based tasks, e.g., plant cover estimation.

# Investigating Neural Network Training on a Feature Level using Conditional Independence

**Author:** Niklas Penzel[1]

**Co-authors:** Christian Reimers[2]; Joachim Denzler[1]; Paul Bodesheim[1]

[1] *Friedrich Schiller University Jena*

[2] *Max Planck Institute for Biogeochemistry*

There are still unresolved questions regarding the changes in learned representations of deep models during the training process. Gaining a better understanding of this process can assist in validating the training. To achieve this goal, previous research has analyzed training in the mutual information plane. We base our analysis on a method founded on Reichenbach's common cause principle. By employing this method, we examine whether the model utilizes information in human-defined features. Given a set of such features, we investigate the changes in relative feature usage throughout the training process. Our analysis includes multiple tasks, e.g., melanoma classification as a real-world application. We discover that as training progresses, models focus on features containing information relevant to the task, resulting in a form of representation compression. Importantly, we also find that the chosen features can differ between training from scratch and fine-tuning a pre-trained network.

# ACQuA: Answering Comparative Questions with Arguments

**Authors:** Alexander Bondarenko[1]; Matthias Hagen[1]

[1] *Friedrich Schiller University Jena*

In the ACQuA project, we develop algorithms to understand and answer comparative information needs like 'Is a cat or a dog a better friend?' by retrieving and combining facts, opinions, and arguments from web-scale resources. Ideally, an answer explains why under what circumstances which comparison alternative should be chosen. Retrieval-based comparative question answering starts with identifying the important constituents: (1) the *objects* that should be compared ('cat' and 'dog' in the above example), (2) the *aspects* that indicate which properties should be emphasized in a comparative answer ('friend'), and (3) *predicates* that guide the direction of the comparison ('better'). When deriving a comparative answer by combining different sources (e.g., different web pages), the following steps can be important: (1) relevance assessment of the individual sources (e.g., a web forum on pets might be more relevant than a page on cat or dog movies), (2) quality assessment and stance detection (e.g., pro 'cat' or pro 'dog') of the retrieved arguments, (3) argument clustering based on the semantic similarity, stance, and quality, (4) re-ranking based on the predicted stance and quality, and (5) answer generation from the final ranking. So far, our fine-tuned RoBERTa-based token classifier (trained and evaluated on 3,500 manually labeled comparative questions) can very reliably identify comparison predicates (almost perfect F1 of 0.98) and objects (F1 of 0.93), while aspect identification falls a bit behind (F1 of 0.80) [1]. Our sentiment-prompted RoBERTa-based stance detector (trained and evaluated on 950 manually labeled answers) still leaves quite some room for improvement (accuracy of 0.63) [1]. For questions that do not contain explicit objects or aspects (e.g., "What pet is best?"), we currently develop approaches that generate clarifying questions and refine the search results based on the feedback (our user study has shown that clarifying comparisons helps [5]). We have also developed "argumentativeness" axioms [2, 4] that help to rerank documents based on (1) the number of argument units (premises and claims identified with our argument mining tool TARGER [7]), (2) the position of query terms in the argument units, (3) (comparative) argument stance, and (4) rhetorical argument quality. Our first findings from participating at several TREC shared tasks and organizing the Touché argument retrieval shared tasks [3] indicate that such argumentativeness facets are promising to improve rankings for argumentative information needs. However, our first results still leave room for further improvements. For instance, formulating new axioms that consider other argumentativeness facets or argument quality dimensions. Finally, based on the aforementioned components (e.g., semantic argument similarity (argument clusters), stance , and quality), we will work on a concise abstractive answer generation / summarization from the "most relevant" arguments in the retrieved web pages. We will adapt the BiLSTM-based abstractive snippet generation framework of Chen et al. [6] to combine different relevant arguments into one concise answer snippet.

**Reference**

[1] Bondarenko, A., Ajjour, Y., Dittmar, V., Homann, N., Braslavski, P., Hagen, M.: Towards understanding and answering comparative questions. In: Proceedings of WSDM 2022. pp. 66–74. ACM (2022)

[2] Bondarenko, A., Fröbe, M., Gohsen, M., Günther, S., Kiesel, J., Schwerter, J., Syed, S., Völske, M., Potthast, M., Stein, B., Hagen, M.: Webis at TREC 2021: Deep Learning, Health Misinformation, and Podcasts tracks. In: Proceedings of TREC 2021. NIST (2021)

[3] Bondarenko, A., Fröbe, M., Kiesel, J., Syed, S., Gurcke, T., Beloucif, M., Panchenko, A., Biemann, C., Stein, B., Wachsmuth, H., Potthast, M., Hagen, M.: Overview of Touché 2022: Argument retrieval. In: Proceedings of the Working Notes of CLEF 2022. CEUR WP, vol. 3180, pp. 2867–2903. CEUR-WS.org (2022)

[4] Bondarenko, A., Fröbe, M., Reimer, J.H., Stein, B., Völske, M., Hagen, M.: Axiomatic retrieval experimentation with ir_axioms. In: Proceedings of SIGIR 2022. pp. 3131–3140. ACM (2022)

[5] Bondarenko, A., Shirshakova, E., Hagen, M.: A user study on clarifying comparative questions. In: Proceedings of CHIIR 2022. pp. 254–258. ACM (2022)

[6] Chen, W., Syed, S., Stein, B., Hagen, M., Potthast, M.: Abstractive snippet generation. In: Proceedings of WWW 2020. pp. 1309–1319. ACM / IW3C2 (2020)

[7] Chernodub, A.N., Oliynyk, O., Heidenreich, P., Bondarenko, A., Hagen, M., Biemann, C., Panchenko, A.: TARGER: Neural argument mining at your fingertips. In: Proceedings of ACL 2019. pp. 195–200. ACL (2019)

# Automated Provision of Data Science Tools to End-Users using Docker & Kubernetes

**Authors:** Julian Hniopek[1]; Nazar Stefaniuk[1]; Thomas Bocklitz[2]

[1] *Friedrich Schiller University Jena & Leibniz Institute of Photonic Technology Jena*

[2] *Leibniz Institute of Photonic Technology, Member of Leibniz Health Technologies, Member of the Leibniz Centre for Photonics in Infection Research (LPI)*

Methods of data science, especially those dealing with machine learning are irreplaceable tools in may fields of science today. Without these methods, many experimental or observational results could not be transferred to useful scientific results. However, in many cases the researchers developing data science methods are not the end-users of them. Especially for todays complex models, such as in the field of deep-learning, highly specialized researchers are necessary to develop and implement appropriate methods for a specific analysis task. This necessitates a quick and easy way to deliver those methods to the domain experts that perform the collection of data and have the necessary knowledge to interpret the results obtained by utilizing data science. Moreover, it poses the challenge that often times these domain experts are not experts in programming or interacting with non-GUI interfaces to run programs, which means that the data science researchers need to provide user-friendly access to their tools. To solve these problems we have developed workflows combining tools for web-based GUI provision and tools for automatic provisioning of these tools to the researchers with none or minimal need for system administrator actions. Using Flask / Django and Shiny it is possible to easily create responsive, web-based user-friendly GUIs as a frontend for access to new algorithms or models tailored to a specific task. Using a Docker based development workflow, data science researchers can use git templates to integrate their algorithms into these GUIs. Using these templates, the resulting applications are automatically built using continuous integration / continuous delivery (CI/CD) pipelines and deployed to a Kubernetes based cluster. These workflows include testing and building and packaging the container, deploying the container to a Kubernetes Cluster using Rancher and ArgoCD as well as setting up appropriate SSL secured networking to the application for easy and secure web-access. Together, this workflow allows to deploy finished data science applications to the end user in a few minutes and facilitates rapid updating and changing of data science methods to adjust to a specific task.
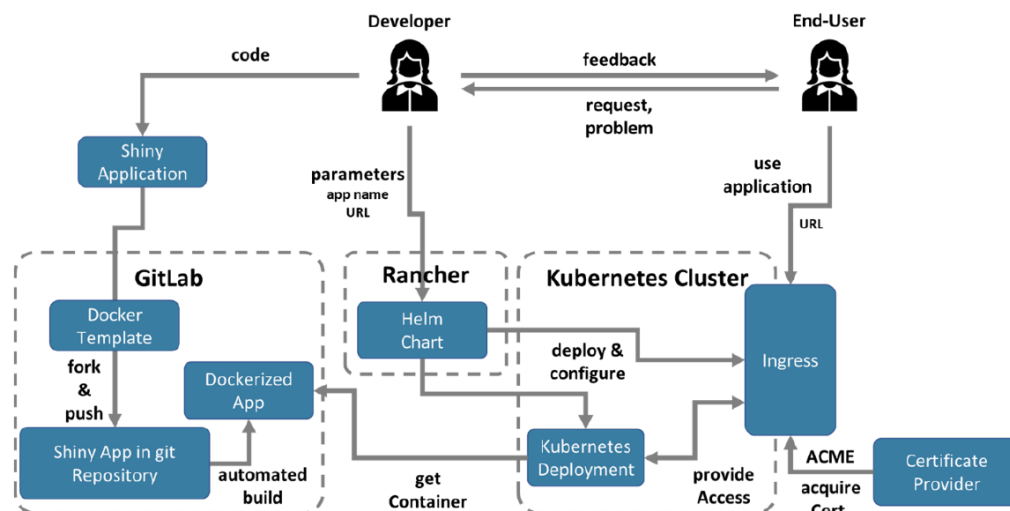


Figure 1: Illustration of the automated data-science application deployment workflow.

# Various Color Transformation Methods

**Authors:** Fatemeh Zahra Darzi; Thomas Bocklitz[1]

[1] *Leibniz Institute of Photonic Technology, Member of Leibniz Health Technologies, Member of the Leibniz Centre for Photonics in Infection Research (LPI)*

**Corresponding Author:** Fatemeh.Darzi@leibniz-ipht.de

Image registration process involves aligning identical shapes or structures in pairs of related images. In order to accomplish this, one image must be transformed into another. The selection of the appropriate transformation model is essential for accurate registration and it should be selected based on the data. In this study, we focus on color transformation methods and apply two different approaches to our data. The first method is based on Reinhard et al. paper[1], which utilizes the Lab color space and the means and standard deviations of each channel to transfer colors from one image to another. However, this method has limitations, such as not preserving the background luminance of the source image. The second method we used is CycleGAN, a type of generative adversarial network (GAN) for image-to-image translation, as proposed by Zhu et al[2]. This approach learns to translate images between source and target domains without paired examples. This approach uses two generators to create images that look real and two discriminators to check if the images are real or fake. However, this approach may not be effective for all tasks, particularly when the source and target domains are significantly different. Another limitation is the time it takes for the model to train. Further research is needed to explore the suitability and limitations of these color transformation methods in different image registration scenarios.
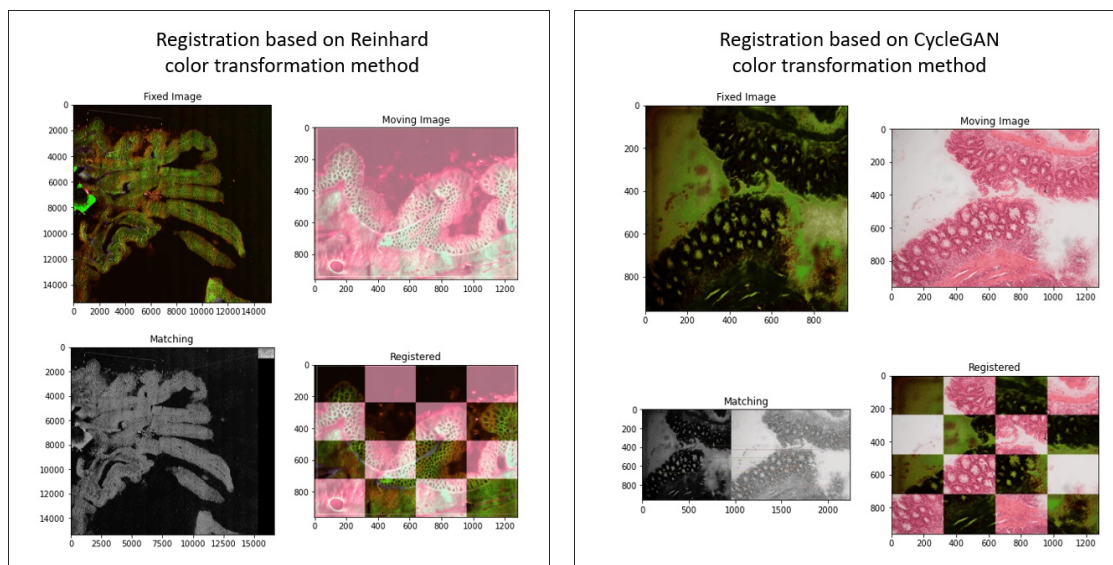


Figure 1: Results of image registration with different color transformation methods.

### References

[1] E. Reinhard, M. Adhikhmin, B. Gooch, and P. Shirley, "Color transfer between images," IEEE Comput. Graph. Appl., vol. 21, no. 5, pp. 34–41, 2001.

[2] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in Proceedings of the IEEE international conference on computer vision, 2017, pp. 2223–2232.

# Epidemiology of genomic surveillance data: Phylodynamic insights into three different phases of the COVID-19 pandemic in Germany

**Authors:** Ariane Weber[1]; Sanni Översti[1]; Denise Kühnert[2]

[1] *Max Planck Institute of Geoanthropology*

[2] *Max Planck Institute of Geoanthropology; Centre for Artificial Intelligence in Public Health Research, Robert Koch Institute*

The dynamics of the COVID-19 pandemic have exhibited different patterns over the course of the first two years. In Germany, early outbreaks that facilitated onward spread in the country were followed by seasonal waves of infection, caused by environmental and behavioural changes. After the introduction of different SARS-CoV-2 variants of concern (VOC), these temporal infection profiles were extended by variant replacement patterns with successive waves caused by Alpha, Delta and Omicron. To track the evolution of SARS-CoV-2 and thus the emergence of new variants, systematic viral genome sequencing was implemented nationally. However, these genomes also hold valuable information on the transmission history of the virus. To extract this information in a phylodynamic birth-death-sampling framework, we use publicly available sequence data and quantify the transmission dynamics of SARS-CoV-2 in Germany during these different time periods. Starting in early 2020, we date the most recent common ancestor of one of the earliest outbreaks in the country, revealing traces of an initial superspreading event quantified through a change in the inferred reproductive number. Continuing in 2021, we infer a longer temporal profile of the reproductive number of SARS-CoV-2 in parts of Bavaria. This highlights the local impact of interventions and seasons on the reproductive number. From SARS-CoV-2 genomes sampled in Berlin in December 2021 we finally estimate the relative increase in the transmission rate of the Omicron compared to the Delta variant in the area, evincing a lower Omicron transmission advantage than in the global context. Taken together, these results provide further insights into and quantifications of the transmission as well as evolutionary dynamics of SARS-CoV-2 over the course of the pandemic in Germany, relying only on genomic sequences and their sampling dates as data source.

# LIBS data analysis

**Authors:** Pegah Dehbozorgi[1,2]; Ludovic Duponchel[3]; Vincent Motto-Ros[4]; Thomas Bocklitz[1,2,5]

[1] *Leibniz Institute of Photonics Technology, Member of Leibniz Health Technologies, Member of the Leibniz Centre for Photonics in Infection Research (LPI), Albert-Einstein-Strasse 9, 07745 Jena, Germany.*

[2] *Institute of Physical Chemistry (IPC) and Abbe Centre of Photonics (ACP), Friedrich Schiller University Jena, Member of the Leibniz Centre for Photonics (LPI), Helmholtzweg 4, 07743 Jena, Germany.*

[3] *Univ. Lille, CNRS, UMR 8516 – LASIRE – Laboratoire de Spectroscopie pour Les Interactions, La R´eactivit´e et L'Environnement, Lille, F-59000, France.*

[4] *Institut Lumière Matière, UMR5306 Université Lyon 1-CNRS, Université de Lyon 69622 Villeurbanne cedex, France.*

[5] *Institute of Computer Science, Faculty of Mathematics, Physics & Computer Science, University of Bayreuth Universitaetsstraße 30, 95447 Bayreuth, Germany.*
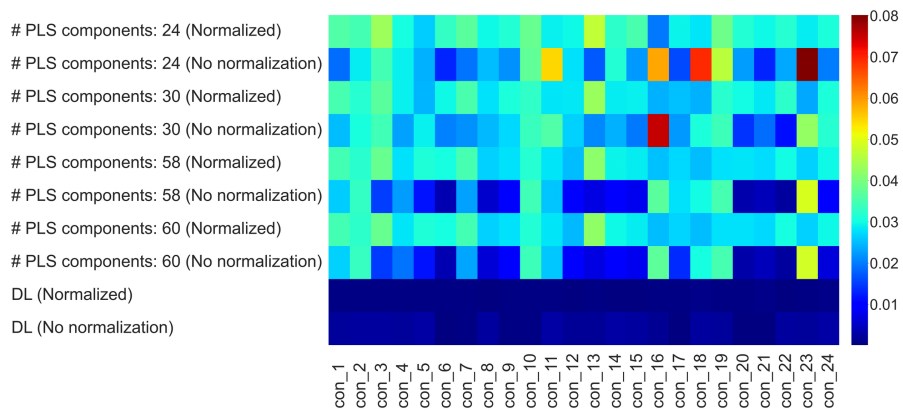
The Laser-Induced Breakdown Spectroscopy (LIBS) technique is widely used to measure the concentration of elements in different types of samples [1, 2]. This study was established to investigate and compare the performance of two approaches, classical regression using Partial Least Square (PLS) and Deep Learning (DL), in predicting the concentration of 24 elements from LIBS spectra. The main challenge for developing predictive models was the variation of electron density and temperature of the plasma, which can completely modify the spectra. Therefore, besides PLS, we tried implementing more advanced tools such as CNNs. The study used the training set of 20000 simulated LIBS spectra and 5000 simulated LIBS spectra as the test set. To develop the models, a pre-processing step was conducted to normalize the data to the (0,1) range. However, the models were also trained and tested with the original data (without normalization) to make the study more comprehensive. For DL, a simple Convolutional Neural Network (CNN) with six convolutional layers was designed. The performance of the models was evaluated based on their stability and accuracy in predicting the concentration of the 24 elements within the test set. Our findings suggest that DL outperformed classical regression in predicting the concentration of presented elements within the simulated test LIBS spectra. The DL model showed greater stability and higher accuracy in predicting concentrations of elements. Overall, this study provides important insight into the application of DL in LIBS analysis as a powerful and stable tool for accurate and reliable elemental analysis.
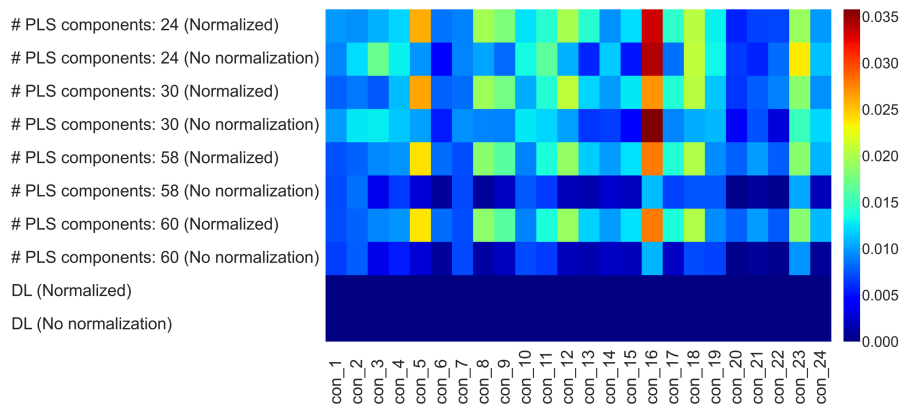
### Acknowledgements

### References

[1] Pasquini, C., et al., Laser induced breakdown spectroscopy. Journal of the Brazilian Chemical Society, 2007. 18: p. 463-512.

[2] Nicolini, O., Libs multivariate analysis with machine learning. 2020.

(a)



(b)

Figure 1: (a) Comparison of the performance of the PLS and DL models to check their stability considering IQR values. (b) Comparison of the performance of PLS models to check their performance regarding the median values.

# Using Advanced Metabolome Annotation Workflow to elucidate Microbial Interactions in Microalgae

**Authors:** Anne-Susann Abel; Mahnoor Zulfiqar[1]

**Co-authors:** Kristy Southysa Syhapanha[1]; Remington X. Poulin; Sassrika Nethmini Costa Warnakulasuriya D.[1]; Georg Alfons Heiner Pohnert[1]; Hans Christoph Steinbeck[1]; Kristian Peters; Maria Sorokina

[1] *Friedrich Schiller University Jena*

**Corresponding Author:** anne-susann.abel@uni-jena.de

Microbial communities reveal complex metabolic interactions which critically influence the ecosystem around them. In this instance, we study the interaction in an aquatic ecosystem between the microalgal species Prymnesium parvum, which produces fish killing toxins, and the diatom Skeletonema marinoi, which causes algal blooms. To decipher the roles of individual microorganisms within their complex community, there is a need to identify the specific metabolites causing the observed effects. Most of these metabolites, however, remain uncharacterized. To address this issue, we use the untargeted metabolomics approach to investigate known and unknown compounds. The two microalgae are grown in a co-culture chamber, separated by a permeable membrane which allows the exchange of exometabolome. The endo- and exometabolome data are extracted from the co-culture samples and analyzed by comparing them to mono-culture samples, which are taken as control. We present here a MS1 pre-processing workflow as addition to the Metabolome Annotation Workflow (MAW) which takes untargeted LC-MS2 spectra as input data. To feed the data generated from LC-MS2 measurements into MAW, pre-processing on MS1 level is crucial and prepares the data for two major tasks: (1) feature annotation and candidate selection by MAW and (2) statistical analysis to find significant features from the complex mix of metabolites. The MS1 pre-processing workflow includes all essential steps for LC-MS1 data analysis, with prior peak picking parameter optimization. The workflow starts with peak detection and retention time correction and yields feature tables, which are then statistically analyzed by PCA, OPLS-DA, variation partitioning and more. Additionally, the MS2 spectra are reconstructed and linked to their MS1 origin for annotation by MAW. Unsupervised statistical analysis of the exometabolome showed significant differences between co-culture and mono-culture samples, suggesting a change in excreted metabolomes for the microorganisms grown in the co-culture experiment. Using variable selection we will identify the features and in turn the metabolites underlying this relationship. The results from these metabolomic analyses can later be integrated into a multi-omics approach and combined with transcriptomics data to give an overview of this experimental microbial co-culture.

# Unsupervised Anomaly Detection for Space Gardening

**Authors:** Ferdinand Rewicki[1,2]; Joachim Denzler[2]; Julia Niebling[1]

[1] *DLR - Deutsches Zentrum für Luft- und Raumfahrt*

[2] *Friedrich Schiller University Jena*

### Abstact

The EDEN Roadmap at DLR aims at building a Bio-regenerative Life Support System (BLSS) for future space missions within the current decade. To ensure the safe and stable operation of the BLSS, the need for automated system monitoring in general and, in particular, robust anomaly detection is apparent. While the abundance of available methods makes it difficult to choose the most appropriate method for a specific application, each method has its strengths in detecting anomalies of different types. The decision becomes even more difficult if annotated data is not available that could be used for model selection. To address this challenge, we compared six unsupervised anomaly detection methods of varying complexity on the UCR anomaly archive benchmark. The goal was to determine whether more complex methods perform better and if certain methods are better suited to specific anomaly types. To validate our findings in the BLSS domain, we applied the best-performing methods to telemetry data collected from the EDEN ISS research greenhouse, which operated from 2018 - 2021 in Antarctica.

### Introduction

Bio-regenerative Life Support Systems (BLSSs) are utilized in habitats to produce plant-based food and close material cycles for respiratory air, water, biomass, and waste. The EDEN NEXT GEN Project, part of the EDEN roadmap at DLR, aims to design a fully integrated ground demonstrator of a BLSS that includes all subsystems. This project builds upon the findings of its predecessor project, EDEN ISS, in which controlled environment agriculture (CEA) technologies were investigated. EDEN ISS was a near closed-loop system built into two 20-foot ISO containers and deployed to the German Antarctic Station Neumayer III in 2017. From 2018 to 2021, crop cultivation, such as lettuces, bell peppers, leafy greens, and various herbs, was studied [6]. To ensure the safe and stable operation of the BLSS, we are investigating methods to mitigate risks regarding system health. Since there is no clear definition of unhealthy system states or sufficient annotated data available for this kind of application, we investigate unsupervised methods for anomaly detection. Choosing the appropriate method from the plethora of available options for a given application is challenging because different methods have different strengths in detecting certain types of anomalies, and the existence of a universal anomaly detection method is a myth [2]. To address this challenge, we compared six unsupervised anomaly detection methods with varying complexities in [5]. Three of these methods are classical machine learning techniques, while the remaining three are based on deep learning. Our central questions in this comparison have been: (1) "Is it worthwhile to sacrifice the interpretability of classical methods for potentially superior performance of deep learning methods?" and (2) "What different types of anomalies are the methods capable of detecting?" We found that the two classical methods, *Maximally Divergent Intervals (MDI)* [1] and *MERLIN* [4], not only performed best, but they also seemed to complement each other in terms of the detected anomaly types [5]. However, as MERLIN suffered from high runtimes, we switched to an improved method for discord discovery called *Discord Aware Matrix Profile (DAMP)* [3]. To validate the results from [5] in the BLSS domain, we are applying MDI and DAMP to a telemetry dataset collected at the EDEN ISS research greenhouse.

### Methods

MDI [1] is a density-based method for offline anomaly detection in uni- or multivariate, spatiotemporal data. Given a multivariate time series $\mathcal{T}$, MDI detects anomalous subsequences by comparing the probability density $p_S$ of a subsequence $S \subseteq \mathcal{T}$ to the density $p_{\Omega(S)}$ of the remaining part of the times series $\Omega(S) := \mathcal{T} \setminus S$ for all subsequences. For more details on MDI, please refer to [1] and [5].

DAMP [3] is a method for offline and online anomaly detection based on discord discovery: Given a subsequence $S$ with length $L$ starting at timestamp $p$, a matching subsequence $M$ starting at timestamp $q$ is called a non-self match to $S$ if $|p - q| \geq L$ [4]. The discord $\tilde{S}$ of a time series $\mathcal{T}$ is

defined as the subsequence with the largest distance $d(\tilde{S}, M_{\tilde{S}})$ from its nearest non-self match $M_{\tilde{S}}$, where $d(\cdot, \cdot)$ is the z-normalized (zero mean and unit variance) Euclidean distance. Advantages of DAMP compared to MERLIN are, that DAMP can be applied effectively online and to multivariate data is well. For details on DAMP, please refer to [3].

**Data**

To validate the findings from [5] in the BLSS domain, we use a subset of the data, collected in the EDEN ISS research greenhouse. The dataset consists of eight time series of sensor readings for carbon dioxide ($CO^2$), relative humidity (RH), photosynthetic active radiation (PAR) and temperature (T) for the year 2020. These variables have been measured at two different places within the greenhouse and belong to the Atmosphere Management Subsystem of EDEN ISS. Each time series has a sampling rate of one data point every 5 minutes ($0.00\bar{3}$ Hz) and a total length of 105408 data points.

**Preliminary Results**

As the discord of the time series is a single subsequence, employ a sliding window approach. Both methods are applied iteratively to a 30-day window by shifting it by one day on each iteration. The score for the newly analyzed day are appended to the score that has been already obtained. The normalized anomaly scores are classified using a threshold of $0.2$. Figure 1 displays the results for temperature readings T1 and T2, with the time series in blue and orange, the detected anomalies by MDI and DAMP highlighted in red and green respectively, and the obtained anomaly scores shown in with the same color coding in the plots below.
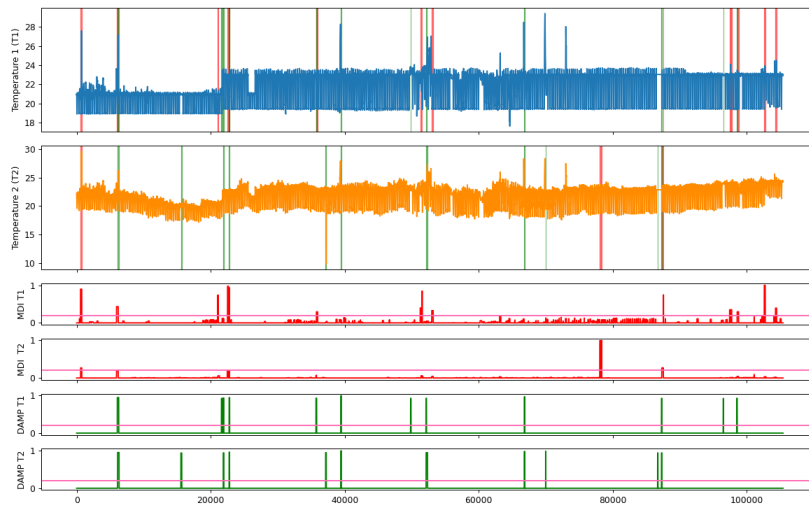


Figure 1:   Sensor readings for T1 (blue) and T2 (orange), MDI anomaly score for T1 and T2 (red) and discords of T1 and T2 as found by DAMP (green).

As there is no ground truth available for anomalies in the EDEN ISS telemetry data, we assess the performance of MDI and DAMP qualitatively. The results confirm that MDI and DAMP identify different types of anomalies. While both methods identify *outlier* anomalies, they do not detect the same instances. Moreover, DAMP identifies *missing drop* anomalies, which manifest as gaps in Figure 1, similar to the findings in [5]. DAMP successfully identifies a *change point* anomaly between time points 20600 and 21753, which is not detected by the MDI method. On the other hand, MDI detects a subtle *local drop* anomaly at time point 21120 that is not identified by DAMP. These results emphasize the usefulness of utilizing both methods in conjunction with each other for effective anomaly detection.

The results for the other variables are similar to those of the temperature readings. MDI and DAMP flag different subsequences as anomalous, which appear reasonable upon visual inspection. However, the dominant pattern of the time series is not as evident as that of the temperature readings. Therefore, the correctness of the detected anomalies should be verified by domain experts.

**Conclusion & Outlook**

Our recent benchmark in [5] indicated that combining MDI with a discord discovery-based anomaly detection method can detect a wide range of different anomalies. The analysis of telemetry data from the EDEN ISS research greenhouse confirms this finding. To ensure a more rigorous evaluation of the results, we will collaborate with BLSS domain experts and obtain their feedback on our initial findings.

**References**

[1] Björn Barz, Erik Rodner, Yanira Guanche Garcia, Joachim Denzler. "Detecting Regions of Maximal Divergence for Spatio-Temporal Anomaly Detection." IEEE Transactions on Pattern Analysis and Machine Intelligence, 2018.

[2] Nikolay Laptev, Saeed Amizadeh, and Ian Flint. Generic and Scalable Framework for Automated Time-series Anomaly Detection. In Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '15). 2015.

[3] Yue Lu, Renjie Wu, Abdullah Mueen, Maria A. Zuluaga, and Eamonn Keogh. DAMP: accurate time series anomaly detection on trillions of datapoints and ultra-fast arriving data streams. Data Min. Knowl. Discov. 2023.

[4] Takaaki Nakamura, Makoto Imamura, Ryan Mercer and Eamonn Keogh. "MERLIN: Parameter-Free Discovery of Arbitrary Length Anomalies in Massive Time Series Archives." IEEE International Conference on Data Mining, 2020.

[5] Ferdinand Rewicki, Joachim Denzler, and Julia Niebling. Is It Worth It? Comparing Six Deep and Classical Methods for Unsupervised Anomaly Detection in Time Series. Applied Sciences 13. 2023.

[6] Zabel, Paul, et al. "Biomass production of the EDEN ISS space greenhouse in Antarctica during the 2018 experiment phase." Frontiers in plant science 11. 2020.

# Fluorescence Lifetime Imaging Microscopy (FLIM) Data Analysis by Inverse Modelling

**Authors:** Mou Adhikari[1]; Thomas Bocklitz

[1] *Friedrich Schiller University Jena*

**Corresponding Author:** mou.adhikari@uni-jena.de

Fluorescence lifetime imaging microscopy (FLIM) provides important information and high-quality images about inter-cellular activity, metabolic state, cellular morphology, etc. [1]. It is a sophisticated imaging approach that relies on the complex curve fitting method by extraction of lifetime parameters. The 'fit-free' deep learning (DL) based lifetime estimation method, which serves as an inverse modelling tool, is the major emphasis of this research. The DL training has been done in two steps: autoencoder and convolutional neural network (CNN). We have carried out our experiments with three datasets to train the autoencoder: (1) noisy data as input and denoised data as output (2) noisy data as input and denoised data as output (without convolution with system response function/IRF) (3) noisy data as input and noisy data as output. After training all, we used the bottleneck features from all three trained autoencoders and used their bottleneck features as input to a CNN to predict lifetime parameters. The last step is the performance analysis of the trained DL model by comparing it with 'FLIMview'[2]. In this study, we also showed our denoising model stability based on different system response functions/ IRF and noise levels. Here, we can see the model performance is quite stable, which represents by the mean square error (MSE) and it is low for all combinations.
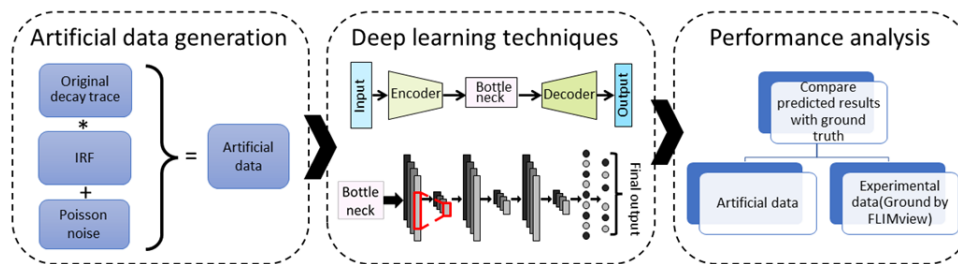


Figure 1: The workflow shows FLIM data analysis by inverse modeling pipeline. In the 1st step, artificial data has been created. The second step is deep learning techniques and training. It has been done in 2 steps: 1. Autoencoder training 2. CNN training. The last step is to test the deep learning model with artificial data and experimental data.

### Acknowledgements

### References

[1] Smith JT, Yao R, Sinsuebphon N, Rudkouskaya A, Un N, Mazurkiewicz J, Barroso M, Yan P, Intes X. Fast fit-free analysis of fluorescence lifetime imaging via deep learning. Proc Natl Acad Sci U S A. 2019 Nov 26;116(48):24019-24030. doi: 10.1073/pnas.1912707116. Epub 2019 Nov 12. PMID: 31719196; PMCID: PMC6883809.

[2] Carrasco Kind M, Zurauskas M, Alex A et al. flimview : A software framework to handle, visualize and analyze FLIM data [version 1; peer review: 1 approved, 1 approved with reservations]. F1000Research 2020, 9:574 (https://doi.org/10.12688/f1000research.24006.1)

# Simplification Models of Scientific and Medical Texts in English and Ukrainian

**Author:** Olha Kanishcheva[1]

[1] *Friedrich Schiller University Jena*

The language used in scientific and medical texts can be difficult for non-experts to understand. In recent years, there has been a growing interest in developing simplification models to make these texts more accessible to a wider audience. In this presentation, we will explore simplification models of scientific and medical texts in both English and Ukrainian. We will discuss the challenges of simplifying such texts and the various approaches that have been used, such as lexical simplification, sentence splitting, and discourse-level simplification. Additionally, we will present some of the tools and resources (ASSET and ASSETUKR) available for the simplification of scientific and medical texts. The findings of this study will be useful for researchers, developers, and educators who are interested in making scientific and medical texts more understandable for the general public.

# FAIR Data for Data Science – Contributions from the FUSION Group

**Authors:** Birgitta König-Ries[1]; the FUSION Group

[1] *Heinz Nixdorf Chair for Distributed Information Systems*

Data Science needs data – more precisely it needs FAIR (Findable, Accessible, Interoperable, Reusable)[1], high quality data. Data Engineering approaches for the provision of such data are an important part of data science.

In our work, we develop such approaches mostly in the context of research data. Ongoing activities include the development of BEXIS2, an open source software for research data management, tools for quality control and curation support, approaches for semantic annotation of datasets, creation and linking of knowledge graphs from data, semantic search and provenance management.

In our poster, we will provide an overview of these approaches along the data life cycle.

**References**

[1] Wilkinson, Mark D., et al. "The FAIR Guiding Principles for scientific data management and stewardship." Scientific data 3.1 (2016): 1-9.

# Data and its challenges on the path to end-to-end digitization in public administration - Contributions from three projects of the openDVA working group

**Authors:** Marianne Mauch[1]; Sarah Bachinger[1]; Sirko Schindler[2]; Leila Feddoul[1]; Felicitas Löffler[3]; Samira Babalou[1] Frank Löffler[1]; Marc Bodenstein[1]; Daniela Erhardt[4]; Clemens Alexander Brust[2];
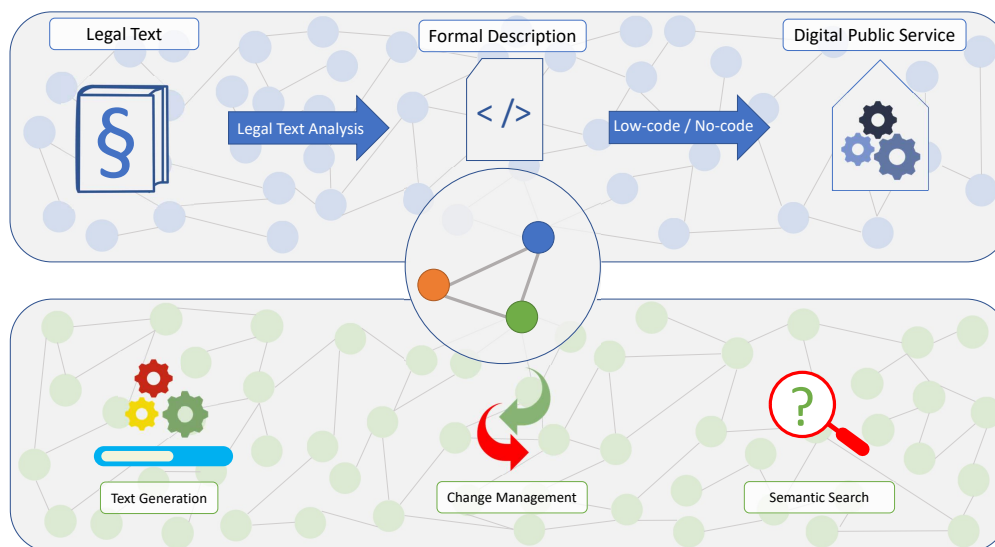
[1] *Friedrich Schiller University Jena*

[2] *German Aerospace Center (DLR)*

[3] *Thuringian Ministry of Finance (TFM)*

[3] *City of Jena*

The implementation of the right to digital access (OZG) in Germany stops at the office door focusing only on the needs of citizens. It does not cover any internal administrative processes and leaves out various stakeholders. For true end-to-end digitization, we need detailed, interoperable descriptions that can be exploited by all interested parties, including small to medium enterprises, decision-makers on all levels, individual administrative staff members, and future citizen developers. They all need a big picture and details on legal regulations, existing standards, and specific requirements. We aim to create such a knowledge base and demonstrate this using a first end-to-end digitized public service. Analyzing structured and unstructured data, for example, in the form of the text of a law addressing a public service, we derive a formal definition of the underlying process and necessary decisions. We enhance this with semantic annotation and link it to available standards. This forms the basis for innovative, new services like a platform for citizen developers to easily create and change fully digitized public services or educational modules that are automatically kept in sync with current developments.



## Acknowledgements

# Interactive Inference - A project at Friedrich Schiller University Jena founded by the Carl-Zeiss-Stiftung

**Authors:** Andreas Kröpelin[1]; Niklas Merk[2]; Vincent Messow[1]; Paul Gerhardt Rump[2]; Philip Schär[1]; Brian Zahoransky[2]; Wolfhart Feldmeier[1]

[1] *University Hospital Jena*

[2] *Friedrich Schiller University Jena*

**Corresponding Author:** inference@uni-jena.de

Learning from experience and making predictions that will guide future actions are at the core of intelligence. These tasks need to embrace uncertainty to avoid the risk of drawing wrong conclusions or making bad decisions. Probability theory offers a framework to represent uncertainty in the form of probabilistic models.

The research training group "Interactive Inference" is an initiative that unites researchers from machine learning and artificial intelligence, algorithm and performance engineering, logic, visualization, and bioinformatics under this paradigm. We work on various topics that range from fundamental algorithms for probabilistic inference such as Markov Chain Monte Carlo via slice sampling, core building blocks of artificial intelligence such as automatic tensor calculus and the relation of probabilistic queries to tensor operations, to applications in structural biology such as the reconstruction and visualization of 3D-structure and molecular dynamics from biomolecular imaging experiments.

**Acknowledgements**