

# xAI in practice: current state of the art, limitations and perspectives

---

Julien Girard-Satabin (CEA LIST): [julien.girard2@cea.fr](mailto:julien.girard2@cea.fr)



# Preliminaries

---

# You and I

## Myself:

1. researcher at CEA on formal methods for software safety and security applied to machine learning;
2. also working on case-based reasoning and out-of-distribution detection in industrial use cases;
3. not a nuclear scientist!

## The audience:

1. fundamental physics practionners;
2. used to computing enormous amount of (structured?) data;

# Definitions

## Explanation

“An explanation is a presentation of (aspects of) the reasoning, functioning and/or behavior of a machine learning model in human-understandable terms” [Nau+23]

“The **belief** (by the trustor) in the ability (of the trustee) to achieve **something**”

# Explanation is a spectrum

Social science have quite a big corpus on what constitutes a good explanation ([Mil19])?

1. *contrastive*: why P instead of Q?
2. *a social process*: A explains P to B
3. *more generic* (cover more facts), *simpler* (quote less causes), and *coherent* (related to previous knowledge) are more easily understood

# Why explaining?



“The software discovered a new fundamental particle with 99% accuracy!”: not enough to convince scientists! What is the causal chain that led to this decision?

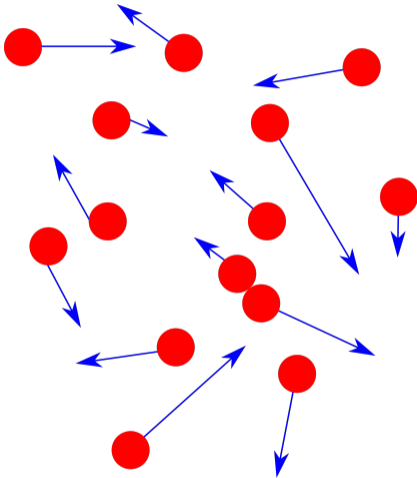
# Why it matters

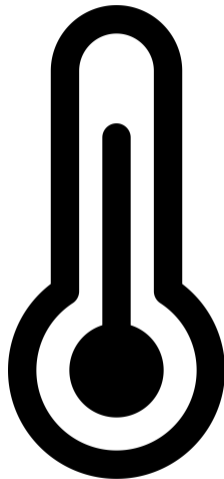
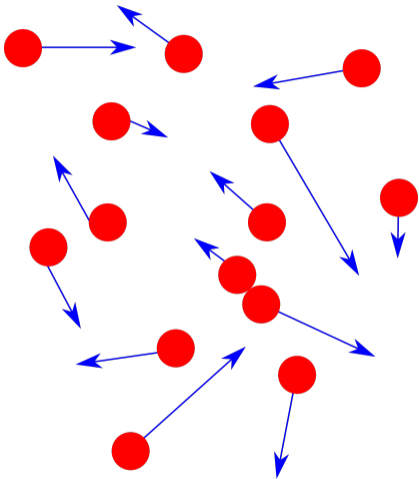
1. debugging and audit
2. refutability
3. compliance with regulation (GDPR article 13.f [SP17])

# About the wording “black-box”

Machine learning is the piling of billions of simple mathematical operations that are atomically well understood







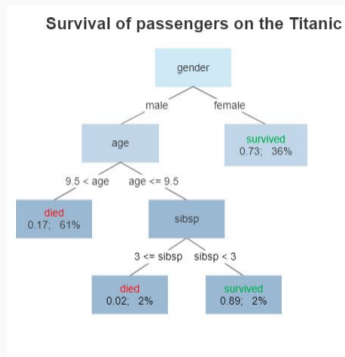
# Post-hoc explanations

---

# Notations

1. samples  $x \in \mathcal{X} \subseteq \mathbb{R}^d$  an input space,  $i^{\text{th}}$  feature  $x_i$
2. an output  $y \in \mathcal{Y} \subseteq \mathbb{R}^p$ , the  $i^{\text{th}}$  feature  $y_i$
3. a program  $f : \mathcal{X} \mapsto \mathcal{Y}$  trained on a  $\mathcal{X}$ 
  - we can usually decompose  $f = h \circ g$
  - in the following,  $h(x)$  is the output of an intermediate layer for neural network
4.  $\nabla_x y$  is the gradient of  $y$  at  $x$

# Decision trees



from Wikipedia [https://en.wikipedia.org/wiki/Decision\\_tree\\_learning/](https://en.wikipedia.org/wiki/Decision_tree_learning/)

Issue: the deeper the tree, the less amenable it is to understand its decision

# Linear regressions

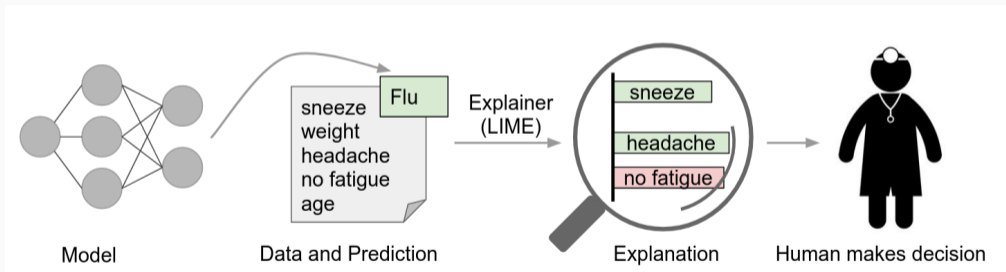
$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \varepsilon$$

A feature will contribute to the decision by its linear coefficient:

$$\beta_k = \frac{y - \sum_{i=1, i \neq k}^{i=n} \beta_i x_i}{x_k}$$

# Under the framework of feature attribution

Basic idea: for a given  $(x, f, y)$ , identify which  $x_i$  was the most useful for the decision



From de [RSG16]

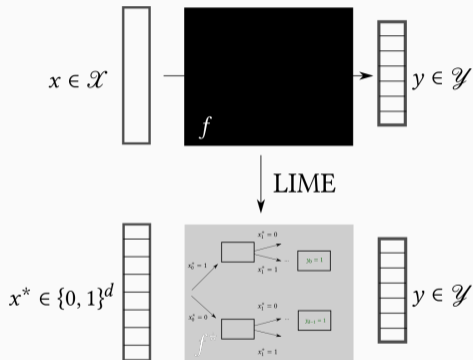
# LIME

## Local Interpretable Model-agnostic Explanations (LIME) [RSG16]:

1. *causal approach*: change  $x_i$  to quantify their impact on  $y$ 
  - if no(sneeze) => no(flu), then sneeze is an important feature
2. once relevant features are identified, train a *surrogate model* that is easier to interpret



# LIME - cont.



The resulting surrogate model only explains *one* prediction

# LIME pros and cons

## Pros:

1. no need for the input data;
2. no need to have access to the program;

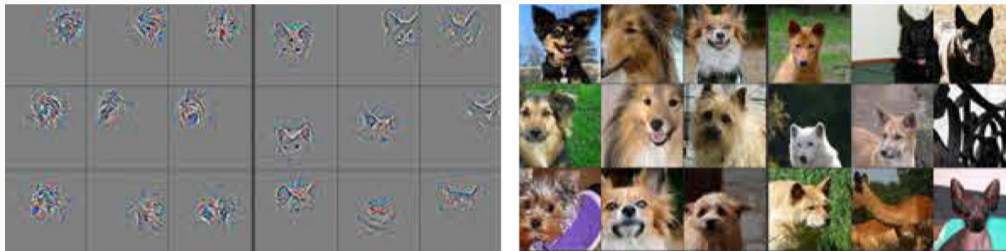
## Cons:

1. training process requires a notion of *neighborhood*, which can be troublesome (images);
2. no validity domain for the surrogate model;

# Derivated approaches: Shapley values

1. identify the mean-shift of each feature contribution SHAP [LL17] (Shapley values) to analyze ensemble models
2. gradually mask parts of the inputs (RISE [PDS18])

# Feature heatmaps



from [ZF14]

basic idea: compute  $\nabla_x y$  and project back on the input space the most important  $x_i$

# GRADCAM, SMOOTHGRAD

GRADCAM [Sel+16; Cha+18] computes  $\nabla_{h(x)} y_i$ , then upsample the resulting point  $\mathcal{X}$

SMOOTHGRAD [Smi+17]  $\nabla_{x^*} y$  where  $x^*$  is a gaussian neighborhood of  $x$

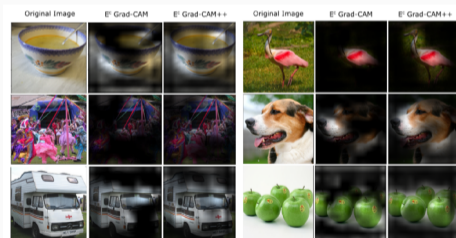


Figure 2: From [Cha+18]

# Integrated gradients





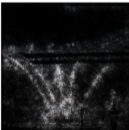
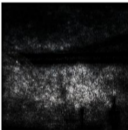



Gradient on the line between  $x$  and a baseline image  $x'$  [STY17]

$$\text{IG}_i = (x_i - x'_i) \int_{\alpha=0}^1 \nabla_{x_i} f(x' + \alpha(x - x')) d\alpha$$

usually computed using Riemann approaches

$$\text{IG}_i \approx (x_i - x'_i) \sum_{k=0}^m \nabla_{x_i} f(x' + \frac{m}{k}(x - x')) * \frac{1}{m}$$

# Integrated gradients

Original image	Top label and score	Integrated gradients	Gradients at image
	Top label: reflex camera Score: 0.993755		
	Top label: fireboat Score: 0.999961		
	Top label: school bus Score: 0.997033		

# Wrapping up: empirical feature attribution approaches

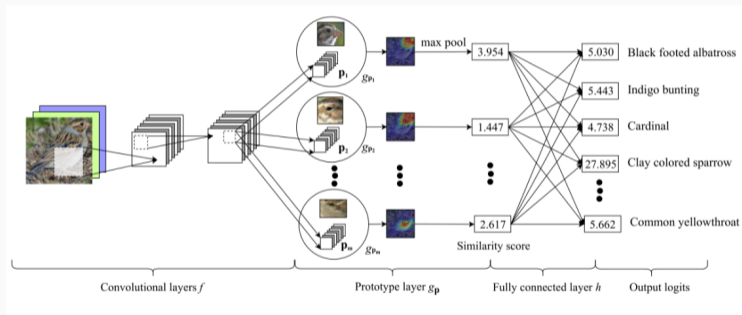
1. usually only require gradient computation access;
2. provide attributions on the input space;
3. heavily rely on the program internal representation;
4. no validity domain;
5. the question of which distance function to use is still open;



# **Explainable by design programs**

---

# Prototype based approaches - ProtoPnet

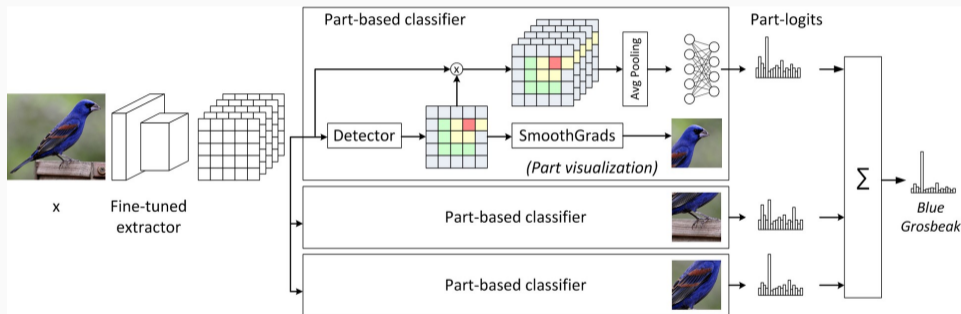


From [Che+19]

# Approches par prototypes - ProtoPnet

1. learn “prototypes” : part of the input set that are used for the prediction;
2. during inference, the various  $h(x)$  are compared to the various prototypes
3. still rely on the hypothesis that “proximity in the latent space equals proximity in the input space”

# Class-wise part detectors [Xu-+23c]



From [Xu-+23c]

# And more...

1. diffusion models [Aug+22]

# Limitations

---

# How to evaluate explanation methods?

Some criterion proposed by [Nau+23] (Co12)

1. *correction*
2. *cohérence* (implementation invariance)
3. *compactness* (size of the explanation)
4. *composability*
5. *controllability*

# How to evaluate explanation metrics?

See [Xu-+23a; Xu-+23d] there is no “one size fits all” metric



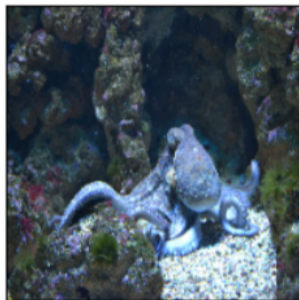
Greyhound (vanilla)



Soup Bowl (vanilla)



Eel (vanilla)



RE 10.9: Images of a dog classified as greyhound, a ramen soup classified as soup bowl, and an octopus classified as eel.

From [Mol22]

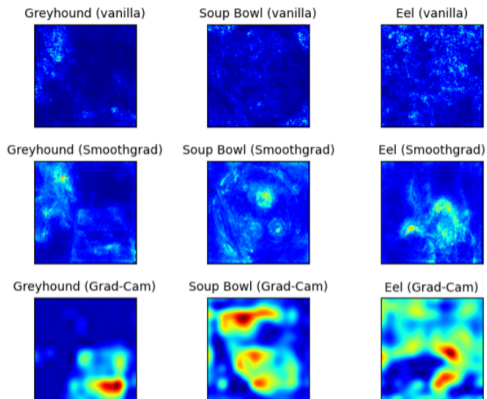


FIGURE 10.10: Pixel attributions or saliency maps for the Vanilla Gradient method, SmoothGrad and Grad-CAM.

From [Mol22]

The network decision is ill-based. *Why* is that? How to fix it?

This explanation does not help to adjust our mental model on the program's behaviour, it is not a good one

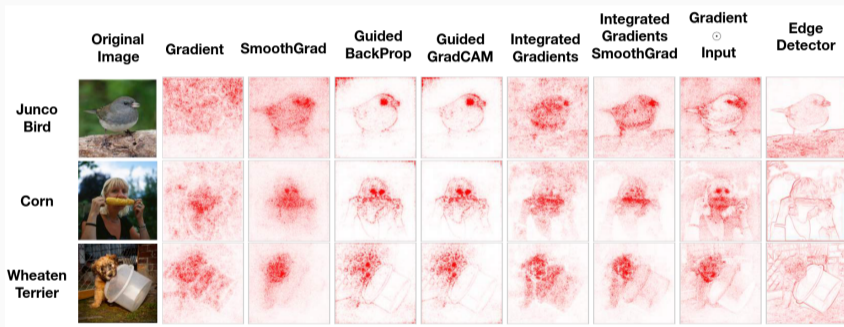
# Nuance

*Extracting a causal chain and displaying it to a person is causal attribution, not (necessarily) an explanation [Mil19].*

Attribution-based approaches are not enough to “fill the holes” for complex programs

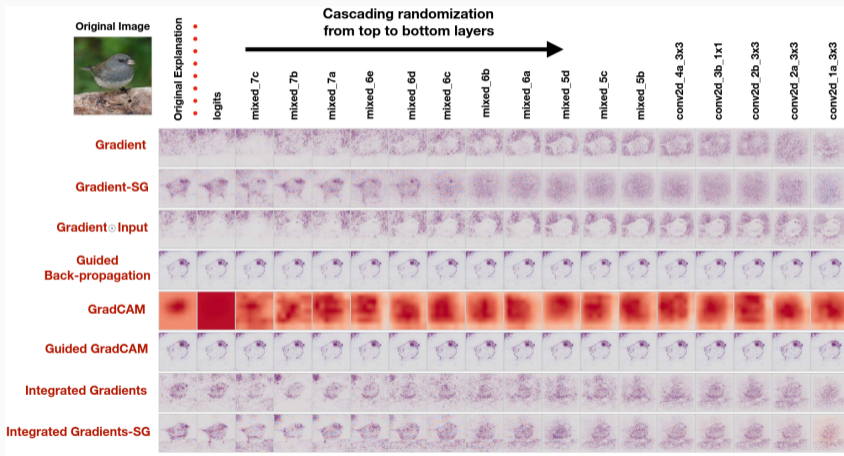
“*How* the decision was taken” and “*Why* the decision was taken” are two different questions

# Feeding our own biases



From [Tom+19]

# Feeding our own biases



From [Tom+19]. The more on the right, the more random the network is.

# Feeding our own biases

## Confirmation bias (Wikitionnaire)

(psychology) A cognitive bias towards confirmation of the hypothesis under study

A “nice” heatmap will confirm that the network works as expected, without being necessarily an accurate description of its inner working

# Explanations can be manipulated [Dom+19]



From [Dom+19]



**The future?**

---

# Open questions

1. validity of feature methods (for a variation on  $f?$  on  $x?$ )
2. how to evaluate explanations and sort evaluation metrics?
3. “social” explanation is yet to happen

# Our work, present and future

1. case-based reasoning [Xu-+23c], out-of-distribution detection [Xu-+23b]
2. explainable by design approaches with a soon-to-come open source library (CABRNET)
3. formal explanation of AI

Open to collaborations!

## References

---

- [Aug+22] Maximilian Augustin, Valentyn Boreiko, Francesco Croce, and Matthias Hein. *Diffusion Visual Counterfactual Explanations*. 2022. DOI: 10.48550/ARXIV.2210.11841. URL: <https://arxiv.org/abs/2210.11841> (cit. on p. 29).

## Bibliography ii

- [Cha+18] Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. “Grad-CAM++: Generalized Gradient-Based Visual Explanations for Deep Convolutional Networks”. In: *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, Mar. 2018. DOI: [10.1109/wacv.2018.00097](https://doi.org/10.1109/wacv.2018.00097). URL: <https://doi.org/10.1109%2Fwacv.2018.00097> (cit. on p. 21).

## Bibliography iii

- [Che+19] Chaofan Chen, Oscar Li, Chaofan Tao, Alina Jade Barnett, Jonathan Su, and Cynthia Rudin. “*This Looks like That: Deep Learning for Interpretable Image Recognition*”. In: *Proceedings of the 33rd International Conference on Neural Information Processing Systems* (2019), pp. 8930–8941 (cit. on p. 26).
- [Dom+19] Ann-Kathrin Dombrowski, Maximillian Alber, Christopher Anders, Marcel Ackermann, Klaus-Robert Müller, and Pan Kessel. “Explanations Can Be Manipulated and Geometry Is to Blame”. In: *Advances in Neural Information Processing Systems*. Vol. 32. Curran Associates, Inc., 2019 (cit. on p. 40).

## Bibliography iv

- [LL17] Scott M. Lundberg and Su-In Lee. “A Unified Approach to Interpreting Model Predictions”. In: *NIPS*. 2017 (cit. on p. 19).
- [Mil19] Tim Miller. “Explanation in Artificial Intelligence: Insights from the Social Sciences”. In: *Artificial Intelligence 267* (2019), pp. 1–38 (cit. on pp. 5, 36).
- [Mol22] Christoph Molnar. *Interpretable Machine Learning. A Guide for Making Black Box Models Explainable*. 2nd ed. 2022. URL: <https://christophm.github.io/interpretable-ml-book> (cit. on pp. 33, 34).

## Bibliography v

- [Nau+23] Meike Nauta, Jan Trienes, Shreyasi Pathak, Elisa Nguyen, Michelle Peters, Yasmin Schmitt, Jörg Schlötterer, Maurice van Keulen, and Christin Seifert. “From Anecdotal Evidence to Quantitative Evaluation Methods: A Systematic Review on Evaluating Explainable AI”. In: *ACM Comput. Surv.* (Feb. 2023). Just Accepted. ISSN: 0360-0300. DOI: 10.1145/3583558. URL: <https://doi.org/10.1145/3583558> (cit. on pp. 4, 31).
- [PDS18] Vitali Petsiuk, Abir Das, and Kate Saenko. “RISE: Randomized Input Sampling for Explanation of Black-box Models”. In: *BMVC*. 2018 (cit. on p. 19).



## Bibliography vi

- [RSG16] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. “Why Should I Trust You?: Explaining the Predictions of Any Classifier”. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (2016) (cit. on pp. 15, 16).
- [Sel+16] Ramprasaath R. Selvaraju, Abhishek Das, Ramakrishna Vedantam, Michael Cogswell, Devi Parikh, and Dhruv Batra. “Grad-CAM: Why did you say that?” In: *ArXiv abs/1611.07450* (2016) (cit. on p. 21).
- [Smi+17] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda B. Viégas, and Martin Wattenberg. “SmoothGrad: removing noise by adding noise”. In: *ArXiv abs/1706.03825* (2017) (cit. on p. 21).

## Bibliography vii

- [SP17] Andrew D Selbst and Julia Powles. “Meaningful information and the right to explanation”. In: *International Data Privacy Law* 7.4 (Dec. 2017), pp. 233–242. ISSN: 2044-3994. DOI: [10.1093/idpl/ipx022](https://doi.org/10.1093/idpl/ipx022). eprint: <https://academic.oup.com/idpl/article-pdf/7/4/233/22923065/ipx022.pdf>. URL: <https://doi.org/10.1093/idpl/ipx022> (cit. on p. 7).

## Bibliography viii

- [STY17] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. “Axiomatic Attribution for Deep Networks”. In: *Proceedings of the 34th International Conference on Machine Learning*. Ed. by Doina Precup and Yee Whye Teh. Vol. 70. Proceedings of Machine Learning Research. PMLR, June 2017, pp. 3319–3328. URL: <https://proceedings.mlr.press/v70/sundararajan17a.html> (cit. on p. 22).
- [Tom+19] Richard J. Tomsett, Daniel Harborne, Supriyo Chakraborty, Prudhvi K. Gurram, and Alun David Preece. “Sanity Checks for Saliency Metrics”. In: *ArXiv abs/1912.01451* (2019) (cit. on pp. 37, 38).

## Bibliography ix

- [Xu-+23a] Romain Xu-Darme, Jenny Benois-Pineau, Romain Giot, Georges Quenot, Zakaria Chihani, Marie-Christine Rousset, and Alexey Zhukov. “On the stability, correctness and plausibility of visual explanation methods based on feature importance”. In: *20th International Conference on Content-based Multimedia Indexing (CBMI 2023)* (2023) (cit. on p. 32).

## Bibliography x

- [Xu-+23b] Romain Xu-Darme, Julien Girard-Satabin, Darryl Hond, Gabriele Incorvaia, and Zakaria Chihani. “Contextualised Out-of-Distribution Detection Using Pattern Identification”. In: *Computer Safety, Reliability, and Security. SAFECOMP 2023 Workshops*. Ed. by Jérémie Guiochet, Stefano Tonetta, Erwin Schoitsch, Matthieu Roy, and Friedemann Bitsch. Cham: Springer Nature Switzerland, 2023, pp. 423–435. ISBN: 978-3-031-40953-0 (cit. on p. 43).

## Bibliography xi

- [Xu-+23c] Romain Xu-Darme, Georges Quénot, Zakaria Chihani, and Marie-Christine Rousset. “PARTICUL: Part Identification with Confidence Measure Using Unsupervised Learning”. In: *Pattern Recognition, Computer Vision, and Image Processing. ICPR 2022 International Workshops and Challenges*. Ed. by Jean-Jacques Rousseau and Bill Kapralos. Cham: Springer Nature Switzerland, 2023, pp. 173–187. ISBN: 978-3-031-37731-0 (cit. on pp. 28, 43).

## Bibliography xii

- [Xu-+23d] Romain Xu-Darme, Georges Quénot, Zakaria Chihani, and Marie-Christine Rousset. “Sanity checks for patch visualisation in prototype-based image classification”. In: *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. 2023, pp. 3691–3696. DOI: 10.1109/CVPRW59228.2023.00377 (cit. on p. 32).

## Bibliography xiii

- [ZF14] Matthew D. Zeiler and Rob Fergus. “Visualizing and Understanding Convolutional Networks”. In: *Computer Vision – ECCV 2014*. Ed. by David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars. Cham: Springer International Publishing, 2014, pp. 818–833. ISBN: 978-3-319-10590-1 (cit. on p. 20).