



Sparse Activations for Interpretable Disease Grading

Presented by DJOUMESSI Kerol

*Kerol Djoumessi¹, Indu Ilanchezian¹, Laura Kühlewein¹, Hanna Faber¹,
Christian Baumgartner¹, Bubacarr Bah², Philipp Berens¹, Lisa M. Koch¹*

¹*University of Tübingen, Germany*

²*African Institute for Mathematical Sciences, South Africa*

July 12, 2023

Motivation

What is interpretability?

Interpretability is the degree to which a human can understand the cause of a decision (Miller 2019).

Motivation

What is interpretability?

Interpretability is the degree to which a human can understand the cause of a decision (Miller 2019).

Why do we need interpretable models?

- Trust AI
- Bias detection

Motivation

What is interpretability?

Interpretability is the degree to which a human can understand the cause of a decision (Miller 2019).

Why do we need interpretable models?

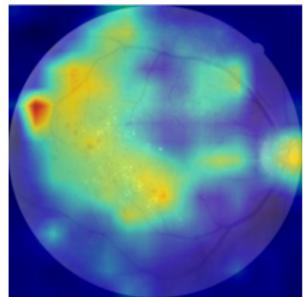
- Trust AI
- Bias detection
- Learning tools
- Policy/legal considerations

Interpretability methods

Why did a model make a specific prediction?

- Attribution maps: where a model looks

Post-hoc explanation

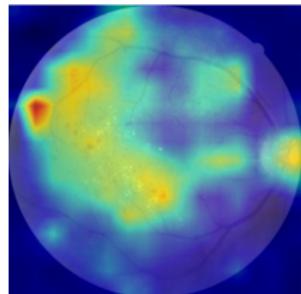


Interpretability methods

Why did a model make a specific prediction?

- Attribution maps: where a model looks

Post-hoc explanation



Interpretable-by-design models

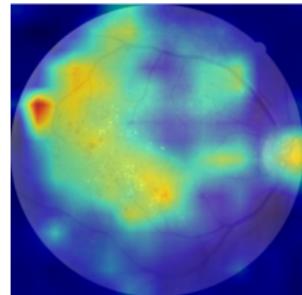
- **Linear models**
- Prototypes-based models

Interpretability methods

Why did a model make a specific prediction?

- Attribution maps: where a model looks

Post-hoc explanation



Interpretable-by-design models

- **Linear models**
- Prototypes-based models

The explanation may look different depending on reason for the explainability.

Attribution maps for medical images: drawbacks

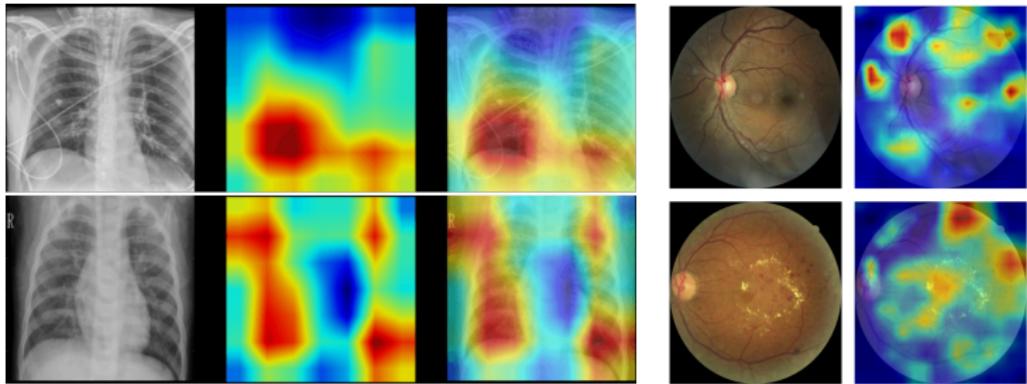


Figure reproduced from Bhowal et al. (2021)

Attribution maps for medical images: drawbacks

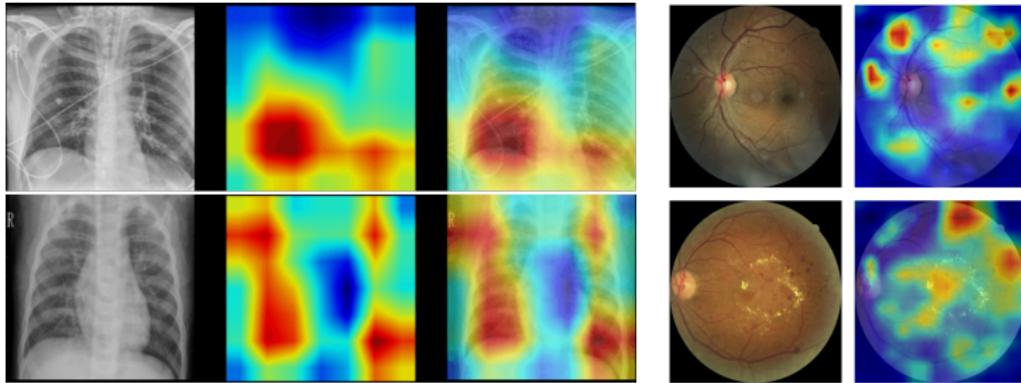


Figure reproduced from Bhowal et al. (2021)

- Highly variable
- Approximate CNNs
- No actionable insights
- Often coarse-grained evidence

Attribution maps for medical images: drawbacks

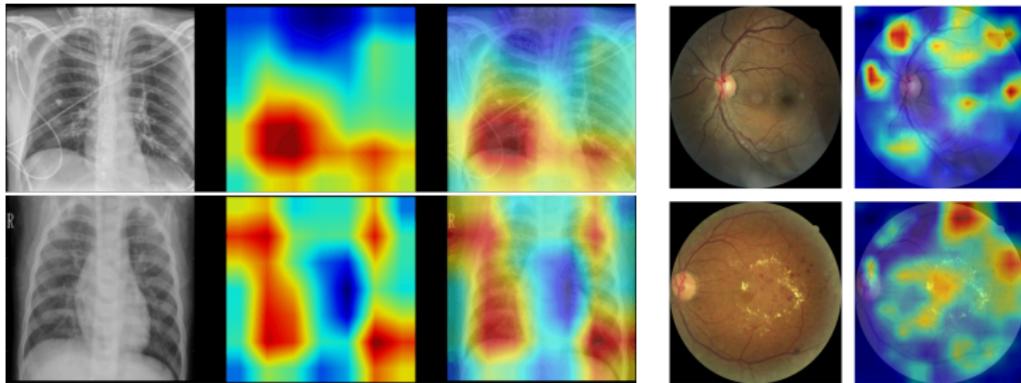


Figure reproduced from Bhowal et al. (2021)

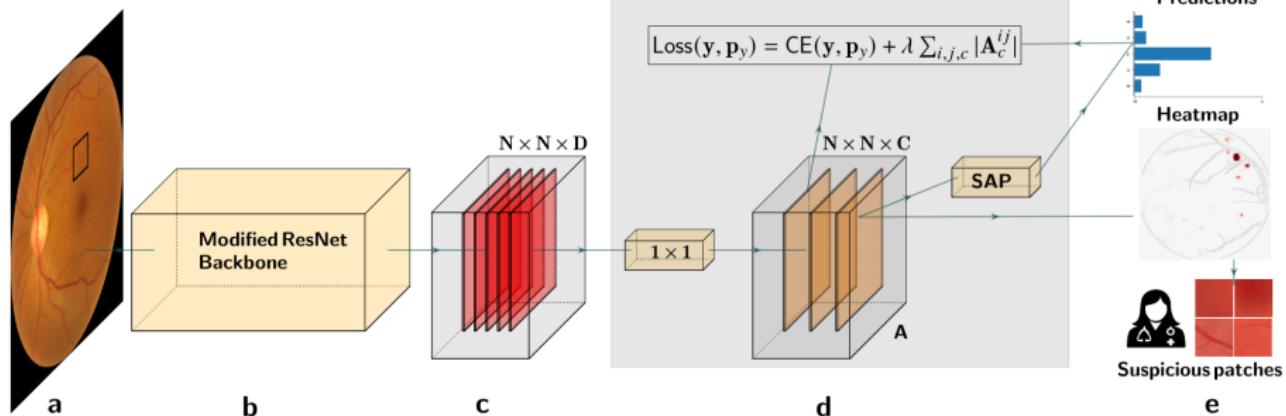
- Highly variable
- Approximate CNNs
- No actionable insights
- Often coarse-grained evidence
- Not inherently interpretable
- Poor performance in Med. data
 - ▶ Reproducibility
 - ▶ Lesion localizations

Contribution: sparse BagNet

We propose an inherently interpretable model that combines the feature extraction capabilities of DNNs with the advantages of sparse linear models in interpretability.

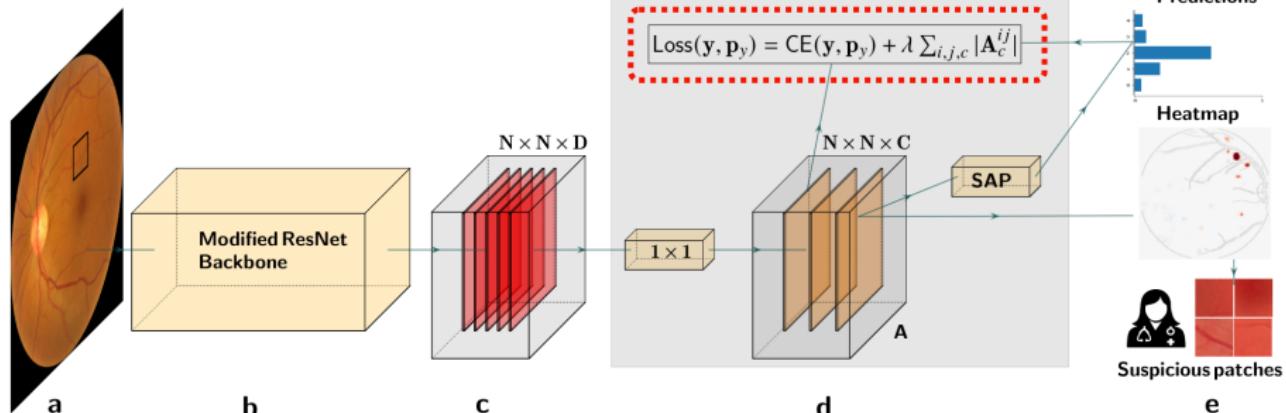
Contribution: sparse BagNet

We propose an inherently interpretable model that combines the feature extraction capabilities of DNNs with the advantages of sparse linear models in interpretability.



Contribution: sparse BagNet

We propose an inherently interpretable model that combines the feature extraction capabilities of DNNs with the advantages of sparse linear models in interpretability.



Application: Diabetic Retinopathy (DR) detection

About DR

- Microvascular complications
- Leading cause of blindness
- Affects 1 in 3 diabetes patients

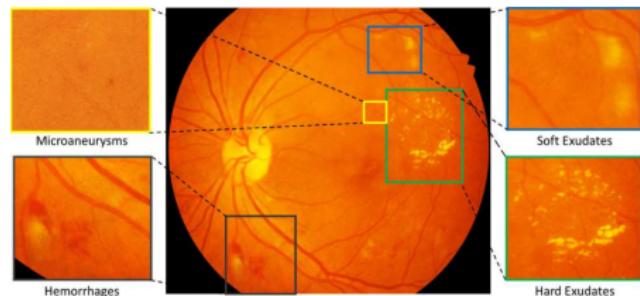
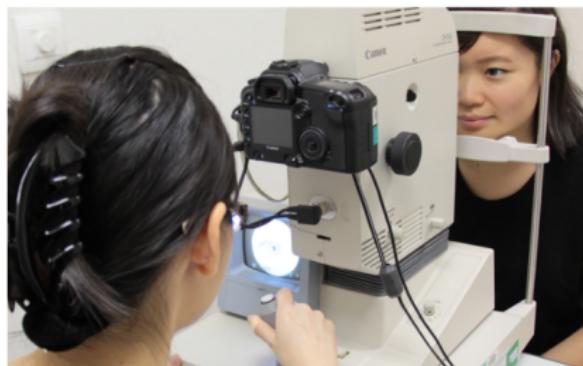
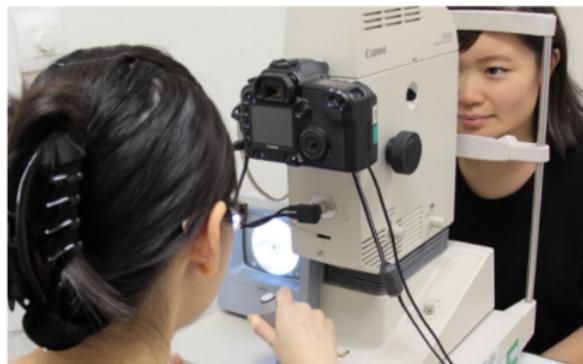


Figure reproduced from Porwal et al. (2018)

Application: Diabetic Retinopathy (DR) detection

About DR

- Microvascular complications
- Leading cause of blindness
- Affects 1 in 3 diabetes patients



Managing DR

- Diagnosis and grading based on retinal fundus images
- Early diagnosis and treatment
- Regular monitoring + screening

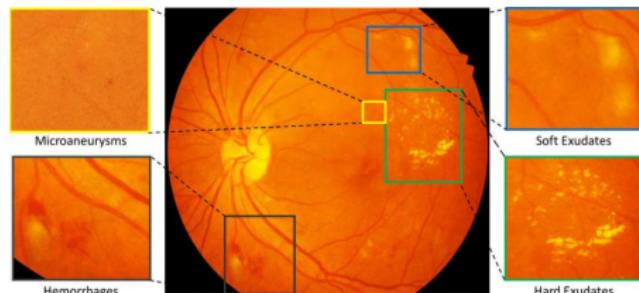


Figure reproduced from Porwal et al. (2018)

Referable DR detection: classification performances

Dataset description:

- Kaggle fundus data set for DR detection
- Referable DR classification task: {0,1} vs {2,3,4}

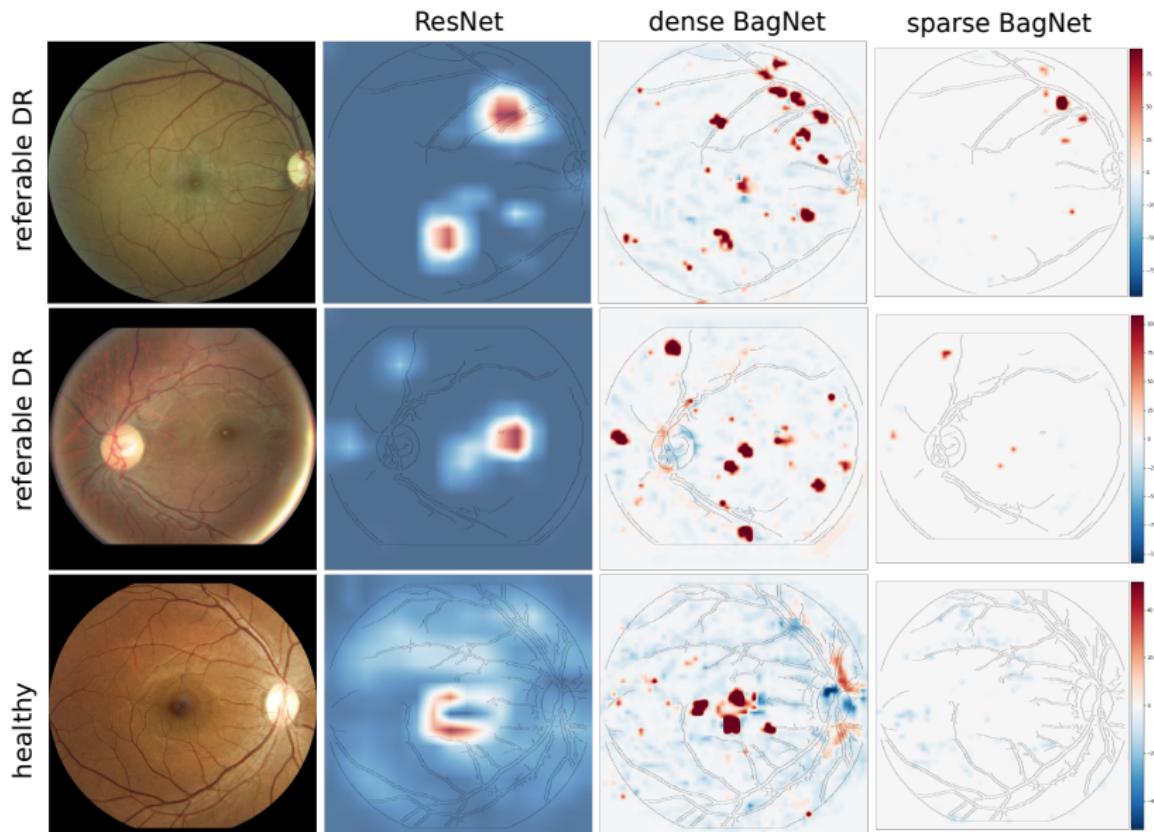
Referable DR detection: classification performances

Dataset description:

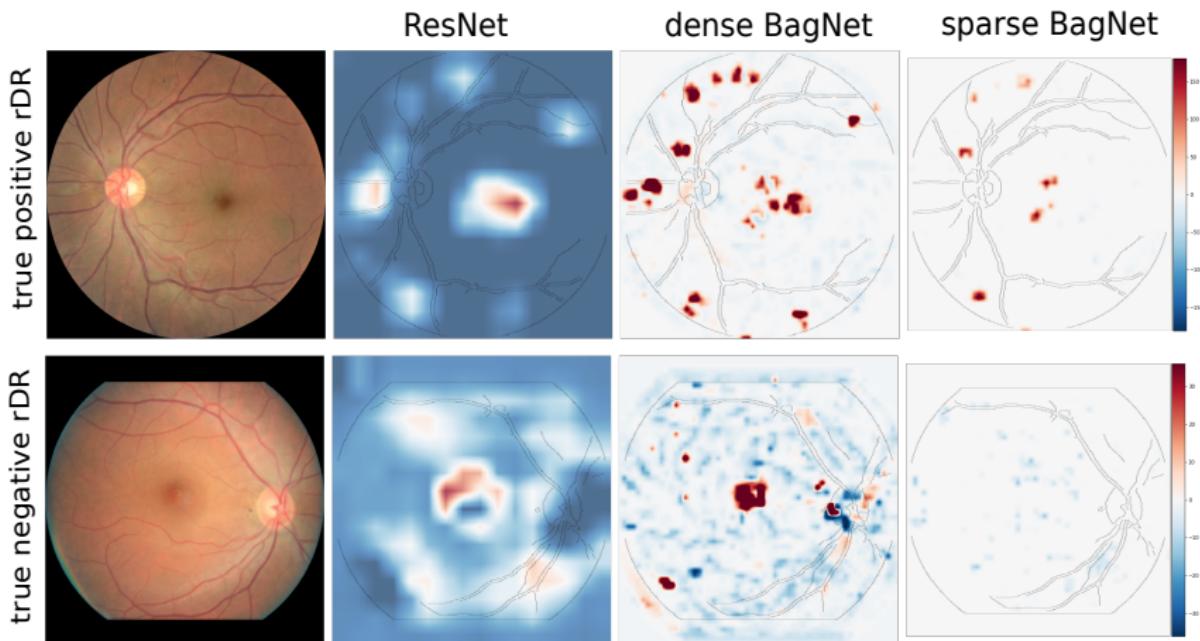
- Kaggle fundus data set for DR detection
- Referable DR classification task: {0,1} vs {2,3,4}

	Accuracy	AUC
ResNet-50	0.942	0.960
Dense BagNet	0.936	0.957
Sparse BagNet	0.928	0.937

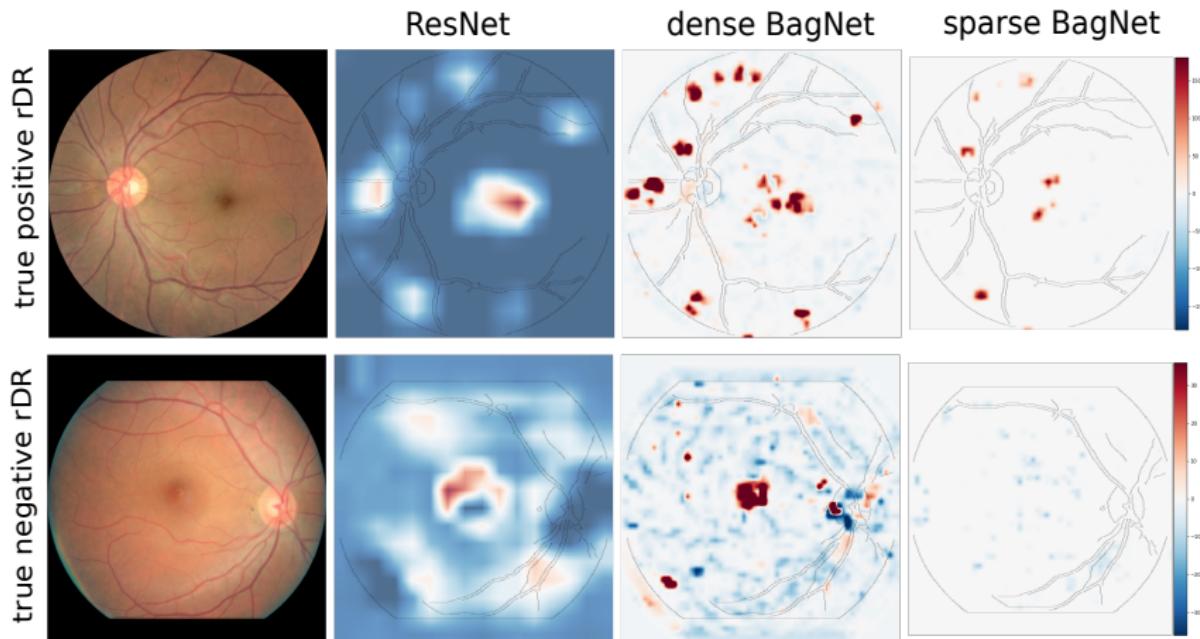
Referable DR detection: qualitative heatmap evaluation



Referable DR detection: quantitative heatmap evaluation



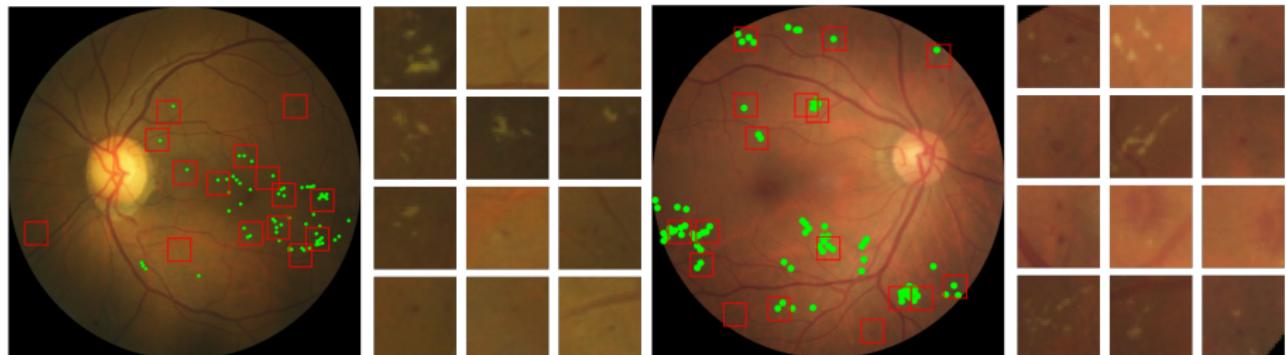
Referable DR detection: quantitative heatmap evaluation



	$\uparrow r_{LG}^-$	$\uparrow r_{LG}^+$
Dense BagNet	0.922 ± 0.03	0.145 ± 0.06
Sparse BagNet	0.991 ± 0.04	0.374 ± 0.33

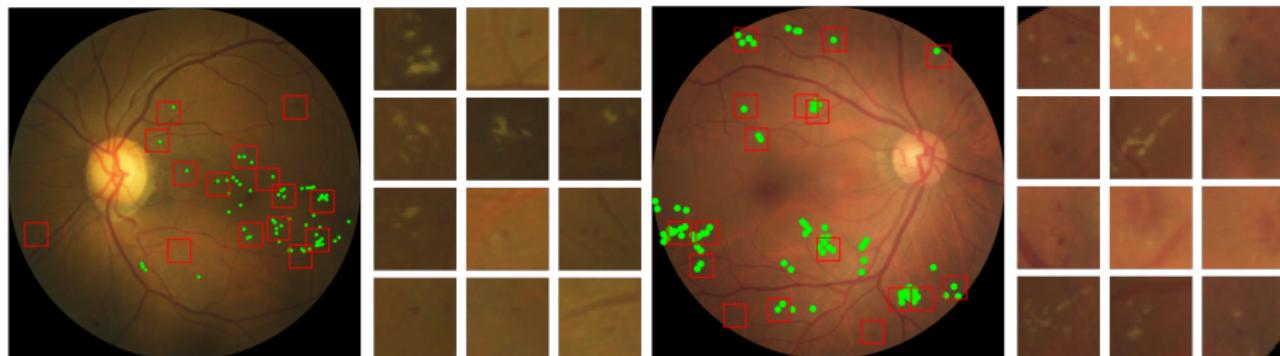
Referable DR detection: quantitative heatmap evaluation

- Lesion annotations (15 images)



Referable DR detection: quantitative heatmap evaluation

- Lesion annotations (15 images)



Precision	
Dense BagNet	0.219 ± 0.1
Sparse BagNet	0.791 ± 0.1

Multiclass DR detection: classification performances

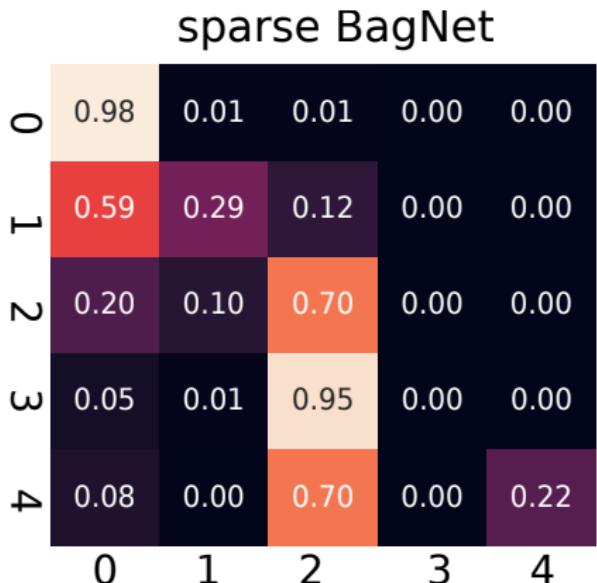
- Multiclass task
- $\{0\}$ vs $\{1\}$ vs $\{2\}$ vs $\{3\}$ vs $\{4\}$

Acc.	
ResNet-50	0.862
Dense BagNet	0.864
Sparse BagNet	0.850

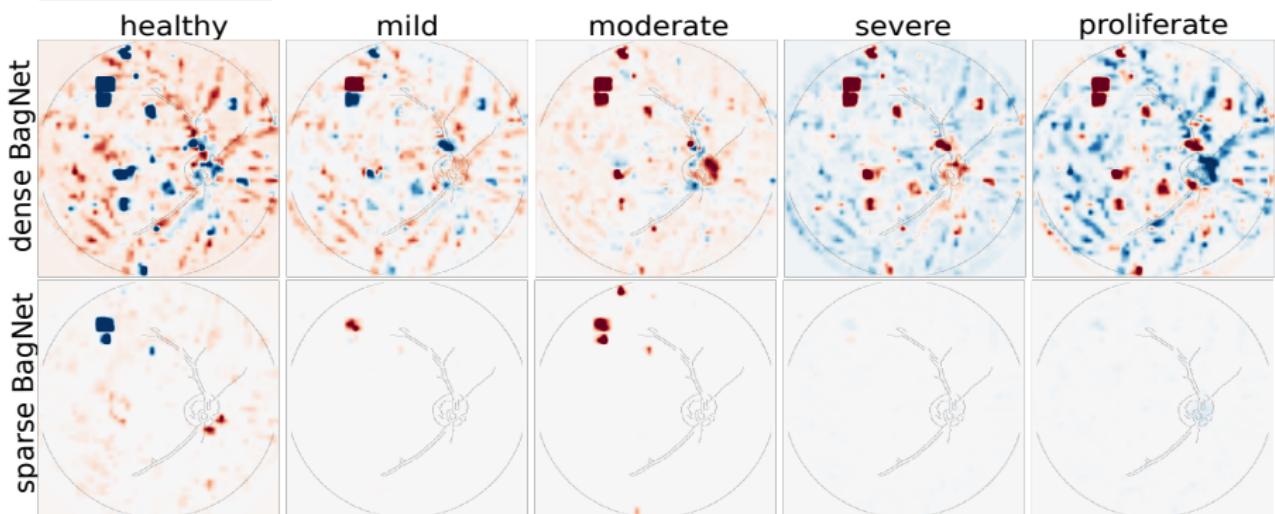
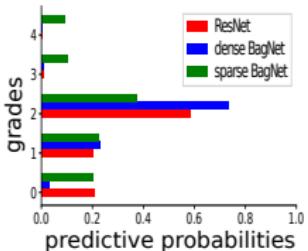
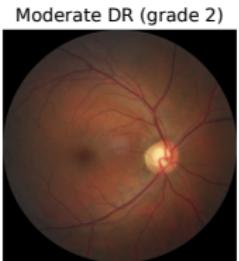
Multiclass DR detection: classification performances

- Multiclass task
- $\{0\}$ vs $\{1\}$ vs $\{2\}$ vs $\{3\}$ vs $\{4\}$

Acc.	
ResNet-50	0.862
Dense BagNet	0.864
Sparse BagNet	0.850



Multiclass DR detection: qualitative heatmap evaluation



Conclusion

Key contributions

- ① BagNet's modifications
- ② Inherently interpretable attribution map
- ③ Quantitative metrics to access the interpretability

Conclusion

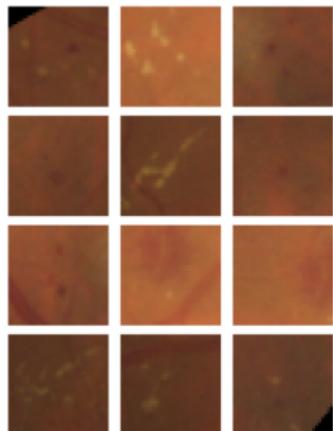
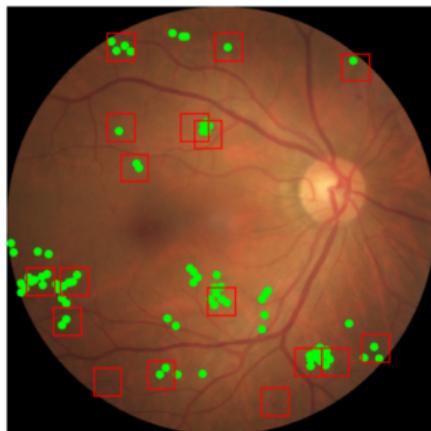
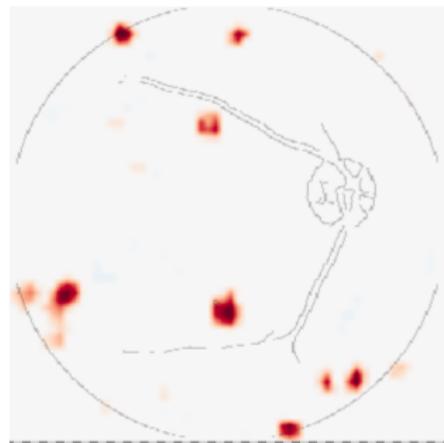
Key contributions

- ① BagNet's modifications
- ② Inherently interpretable attribution map
- ③ Quantitative metrics to access the interpretability

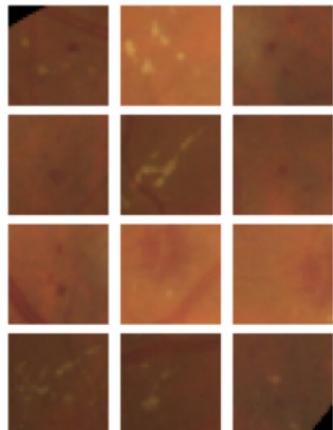
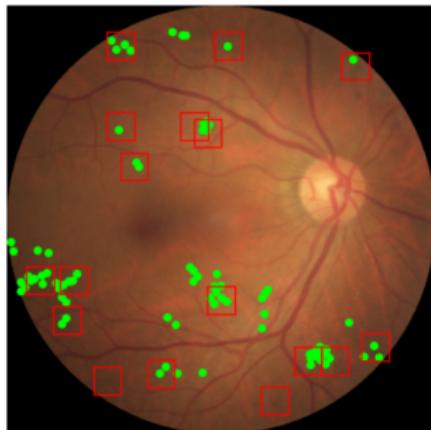
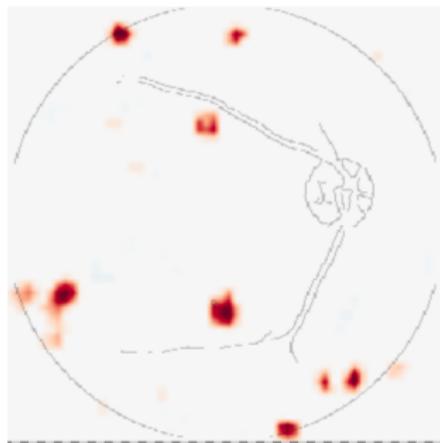
Summary

- Enforcing sparse heatmap improves the interpretability
- Suitable for small lesions distributed over the image
- Comparable SOTA performance on DR detection

Next step



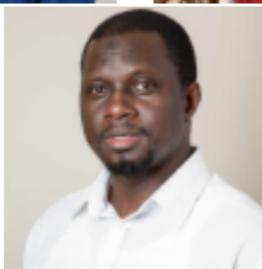
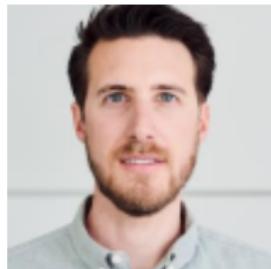
Next step



Clinical study

Does highlighting suspicious areas of the image improve diagnostic performance?

Acknowledgements



Thank you for your attention!

Poster location: **W08**

kerol.djoumessa@uni-tuebingen.de

Dataset

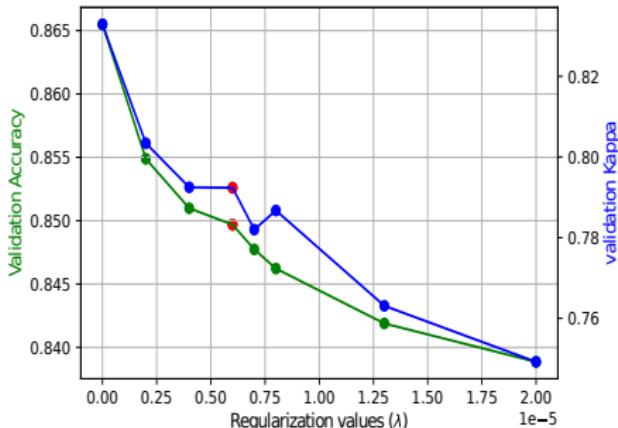
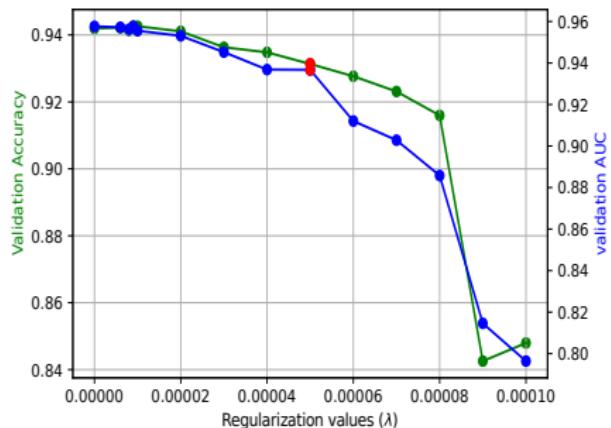
Dataset Description:

- Kaggle dataset: 88,702
- Quality filtering: 45,923
- Referable task: $\{0,1\}$ vs $\{2,3,4\}$
- Multiclass task: $\{0\}$ vs $\{1\}$ vs $\{2\}$ vs $\{3\}$ vs $\{4\}$

Method: Model's architecture

- BagNet encodes the input:
 - ▶ $f_{bag}(X) = A' \in R^{N \times N \times D}$
 - ▶ $A' \xrightarrow{C_1 \times 1} A \in R^{N \times N \times C}$
 - ▶ $A \xrightarrow{\text{SAP}} I_c \in R^{1 \times C}$
- Patch-based model
 - Small receptive field
 - Heatmap = evidence map
 - λ controls the sparsity
 - Suspicious patch visualization

Sparsity hyperparameter λ



Sparse BagNet limitations

ResNet

0	0.98	0.01	0.02	0	0
1	0.56	0.27	0.17	0	0
2	0.17	0.05	0.71	0.07	0
3	0.03	0	0.37	0.6	0
4	0.08	0	0.27	0.26	0.39

dense BagNet

0	0.97	0.01	0.01	0.00	0.00
1	0.53	0.34	0.13	0.00	0.00
2	0.17	0.10	0.68	0.05	0.01
3	0.02	0.01	0.42	0.52	0.03
4	0.03	0.00	0.30	0.25	0.43

sparse BagNet

0	0.98	0.01	0.01	0.00	0.00
1	0.59	0.29	0.12	0.00	0.00
2	0.20	0.10	0.70	0.00	0.00
3	0.05	0.01	0.95	0.00	0.00
4	0.08	0.00	0.70	0.00	0.22