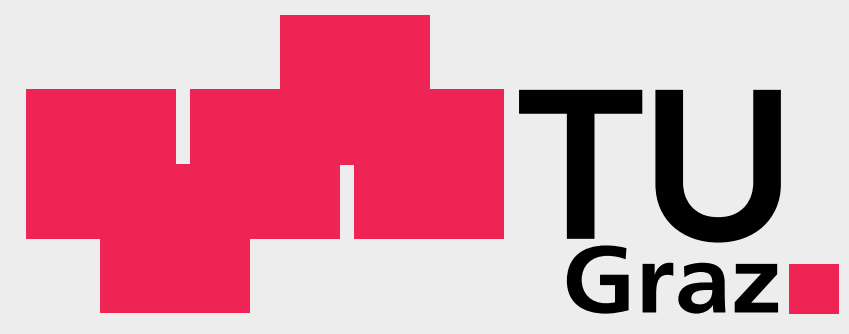


# On Resource-Efficient Bayesian Network Classifiers and Deep Neural Networks



Wolfgang Roth<sup>1</sup>, Günther Schindler<sup>2</sup>, Holger Fröning<sup>2</sup>, and Franz Pernkopf<sup>1</sup>



Der Wissenschaftsfonds.

<sup>1</sup> {roth,pernkopf}@tugraz.at,  
<sup>2</sup> {guenther.schindler,holger.froening}@ziti.uni-heidelberg.de

<sup>1</sup> Signal Processing and Speech Communication Laboratory, Graz University of Technology

<sup>2</sup> Institute of Computer Engineering, Ruprecht Karls University, Heidelberg

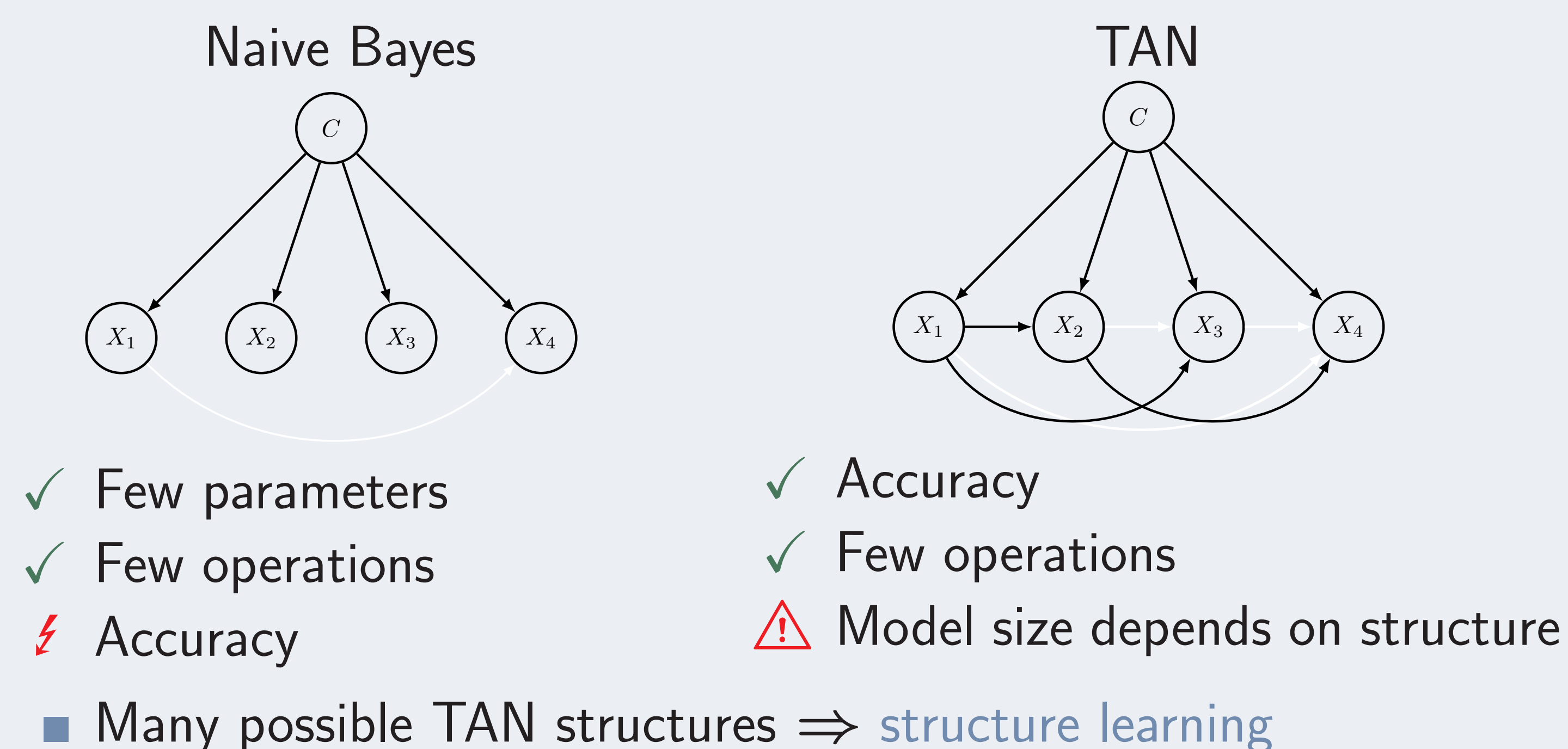
## Summary

- **Reduce model complexity of Bayesian network classifiers**
  - Transfer techniques from deep learning [1] to Bayesian networks
- **Method 1: Model-size-aware TAN structure learning**
  - Apply differentiable TAN structure learning from [2]
    - ✓ Structure learning using backpropagation
  - New extension: Trade off between accuracy and model size
- **Method 2: Quantization-aware Training**
  - Quantize log-probabilities (CPTs) to few bits
  - Apply straight-through gradient estimator
- **Comparing Bayesian networks and deep neural networks**

## Bayesian Network (BN) Classifiers

- BNs define a factorization of a joint distribution via a directed acyclic graph  $\mathcal{G}$  as
$$p(\mathbf{X}) = \prod_{i=1}^D p(\mathbf{X}_i | \text{pa}(\mathbf{X}_i))$$
- This work: Discrete features  $\mathbf{X}_i$ 
  - Parameters: Conditional probability tables (CPTs)  $\theta$
- Classification: Additional distinct class variable  $\mathcal{C}$ 
  - Classify according to  $\arg\max_c \log p(\mathbf{x}, c)$
  - Classification requires only  $(D + 1) \cdot \#(\text{classes})$  additions in log-domain

## Naive Bayes vs. tree-augmented naive Bayes (TAN)



## Model-Size-Aware TAN Structure Learning

- Encode TAN graph  $\mathcal{G}$  using one-hot vectors  $\mathbf{s} = (\mathbf{s}_1, \dots, \mathbf{s}_D)$
  - CPTs for all possible parents:  $\Theta = \{\theta_c\} \cup \{\theta_1, \dots, \theta_D\}$
- $\mathbf{s}_1 = (1)$

$\mathbf{s}_4 = (0, 0, 1, 0)$

$\Theta_1 = \{\theta_{10}\}$

$\Theta_4 = \{\theta_{40}, \theta_{41}, \theta_{42}, \theta_{43}\}$
- $\mathbf{s}_2 = (0, 1)$

$\mathbf{s}_3 = (0, 1, 0)$

$\Theta_2 = \{\theta_{20}, \theta_{21}\}$

$\Theta_3 = \{\theta_{30}, \theta_{31}, \theta_{32}\}$
- Treat parameters  $\Theta$  and structure  $\mathbf{s}$  jointly
$$\log p(\mathbf{X}, \mathbf{C}) = \log p_{\theta_c}(\mathbf{C}) + \sum_{i=1}^D \sum_{j=0}^{i-1} s_{i|j} \log p_{\theta_{ij}}(\mathbf{X}_i | \mathbf{X}_j, \mathbf{C})$$
  - Continuous relaxation of one-hot  $\mathbf{s}_i$  results in probabilities  $\Phi_i$ 
    - Probabilities  $\Phi = (\Phi_1, \dots, \Phi_D)$  define a distribution over TAN graphs  $\mathcal{G}$
    - Optimize expectation with respect to  $\Phi$  and take most probable graph  $\mathcal{G}$ 

$$\mathcal{L}_{\text{SL}}(\Phi, \Theta) = \mathbb{E}_{\mathbf{s} \sim p_{\Phi}} [\mathcal{L}(\Theta, \mathbf{s})] \quad (\text{proposed in [2]})$$
    - $\mathcal{L}_{\text{SL}}$  is differentiable with respect to  $\Phi \Rightarrow$  optimize with backpropagation
  - This work: **New term** penalizes number of parameters

$$\mathcal{L}_{\text{SL}}^{\text{MS}}(\Phi, \Theta) = \underbrace{\mathcal{L}_{\text{SL}}(\Phi, \Theta)}_{\text{proposed in [2]}} + \underbrace{\lambda_{\text{MS}} \mathbb{E}_{\mathbf{s} \sim p_{\Phi}} [\mathcal{L}_{\text{MS}}(\mathbf{s})]}_{\text{trade-off \#model parameters}}$$

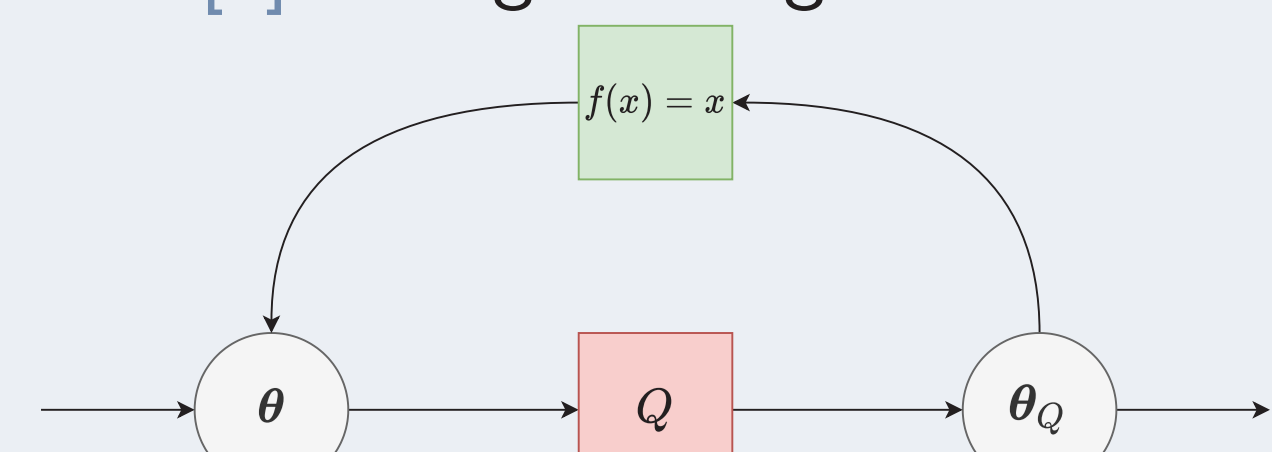
## Quantization-Aware Training

- **BNs:** Quantize log-probabilities to negative **fixed-point values**

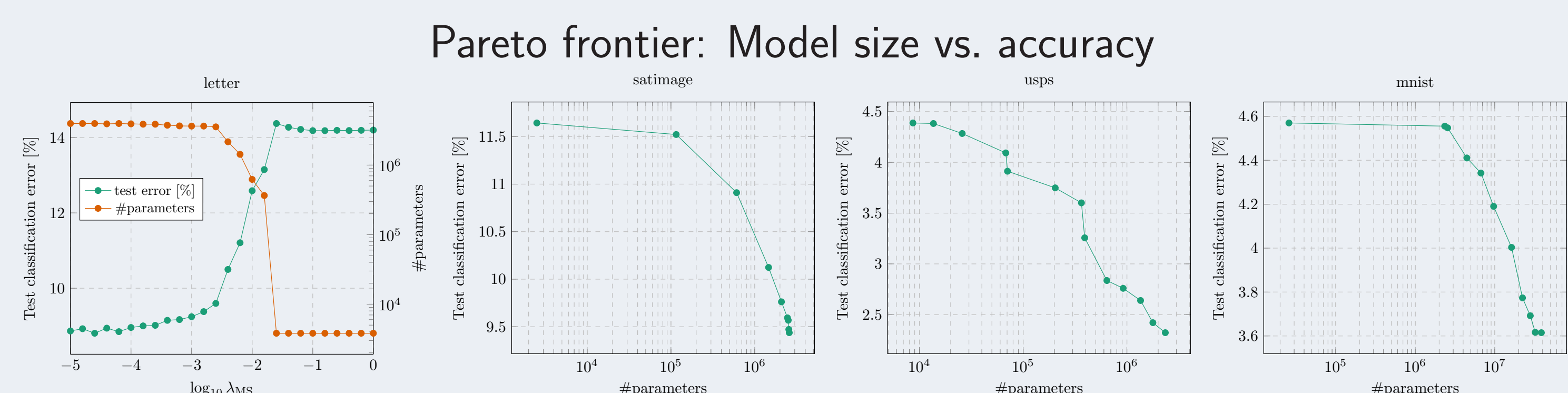
$$Q_{\text{BN}}(\theta) = \text{clip}(\text{round}(\theta \cdot 2^{B_F}) \cdot 2^{-B_F}, -U, 0)$$
- **DNNs:** Quantize weights according to [3]
$$Q_{\text{DNN}}(w) = Q\left(\frac{\text{clip}(w, -1, 1) + 1}{2}; B\right) \cdot 2 - 1$$

$$Q(v; B) = \frac{1}{2^B - 1} \cdot \text{round}((2^B - 1) \cdot v)$$
- The gradient of quantization functions is zero almost everywhere
  - Apply the **straight-through gradient estimator** [4] during training

At backpropagation: Pretend that the gradient of a quantization function is non-zero

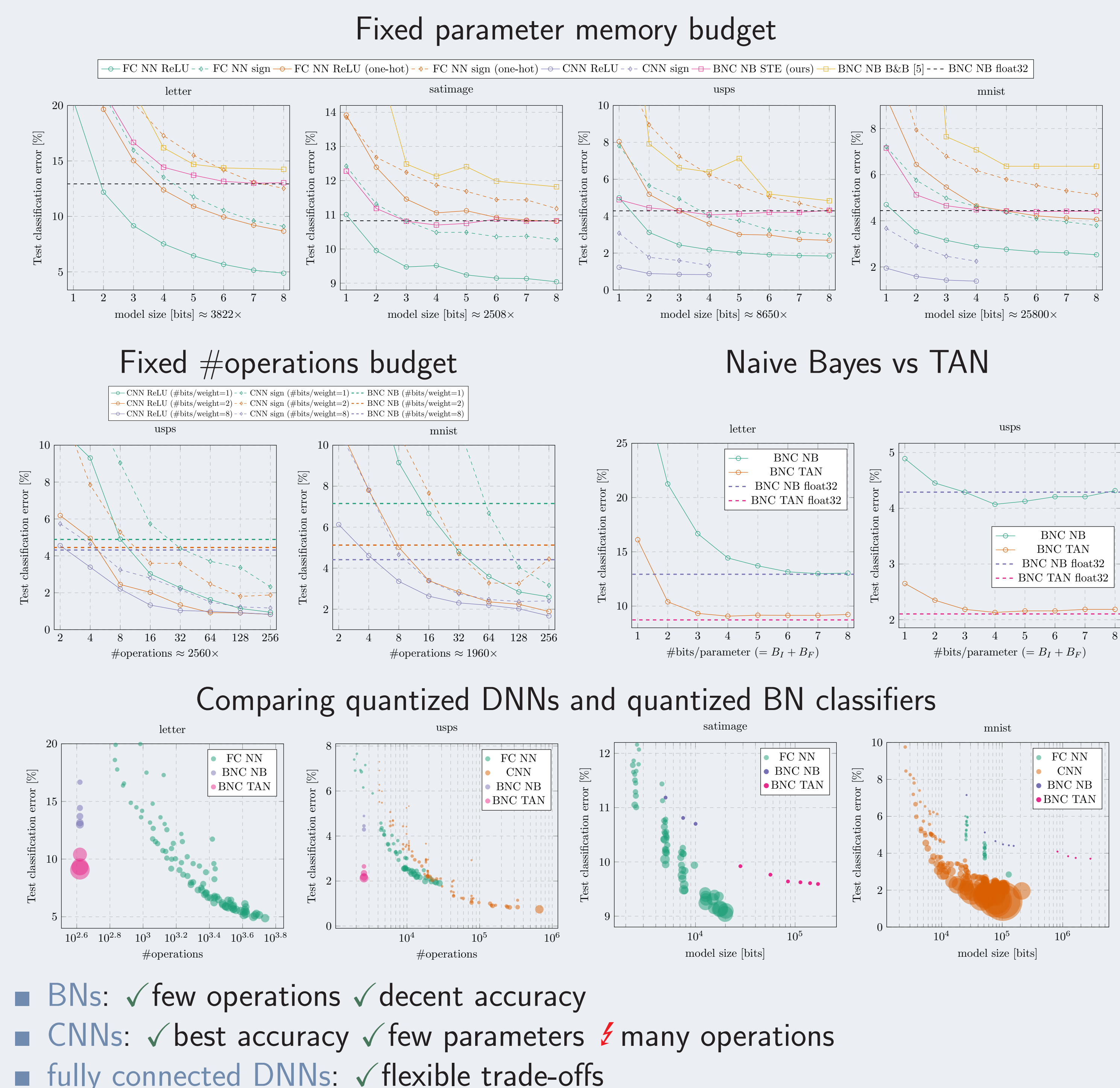


## Experiments: Model-Size-Aware TAN Structure Learning



- ✓ Effective trade-off between accuracy and model size

## Experiments: Quantization-Aware Training



## References

- [1] Roth et al., Resource-Efficient Neural Networks for Embedded Systems, arXiv:2001.03048
- [2] Roth and Pernkopf, Differentiable TAN Structure Learning for Bayesian Network Classifiers, PGM 2020
- [3] Zhou et al., DoReFa-Net: Training Low Bitwidth Convolutional Neural Networks with Low Bitwidth Gradients, arXiv:1606.06160
- [4] Bengio et al., Estimating or Propagating Gradients Through Stochastic Neurons for Conditional Computation, arXiv:1308.3432
- [5] Tschachtschek et al., Integer Bayesian Network Classifiers, ECML 2014