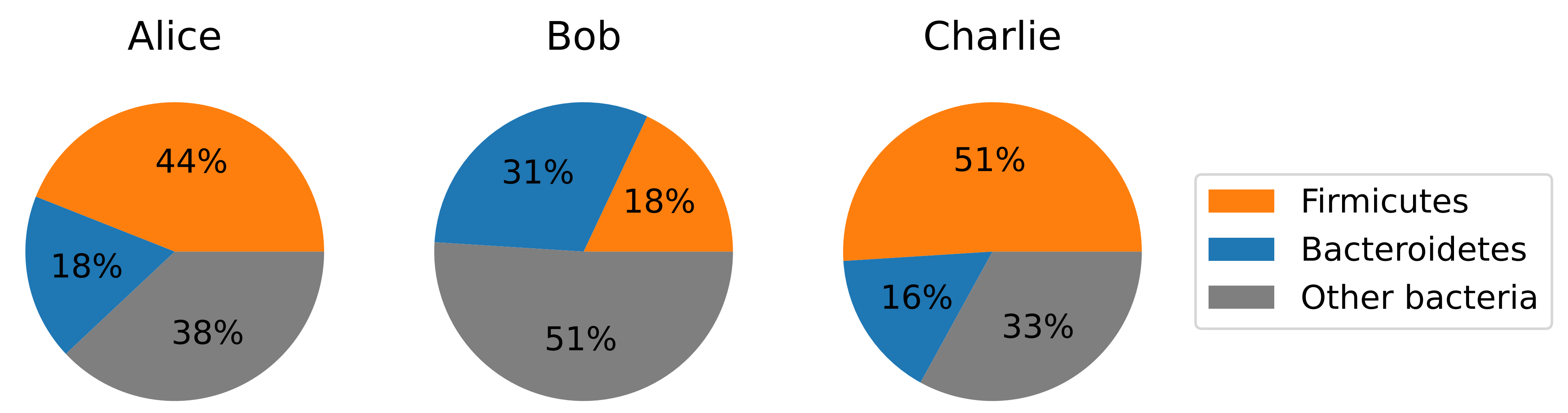


## Compositional Data

Compositional Data (CoDa) describe the parts of a whole:

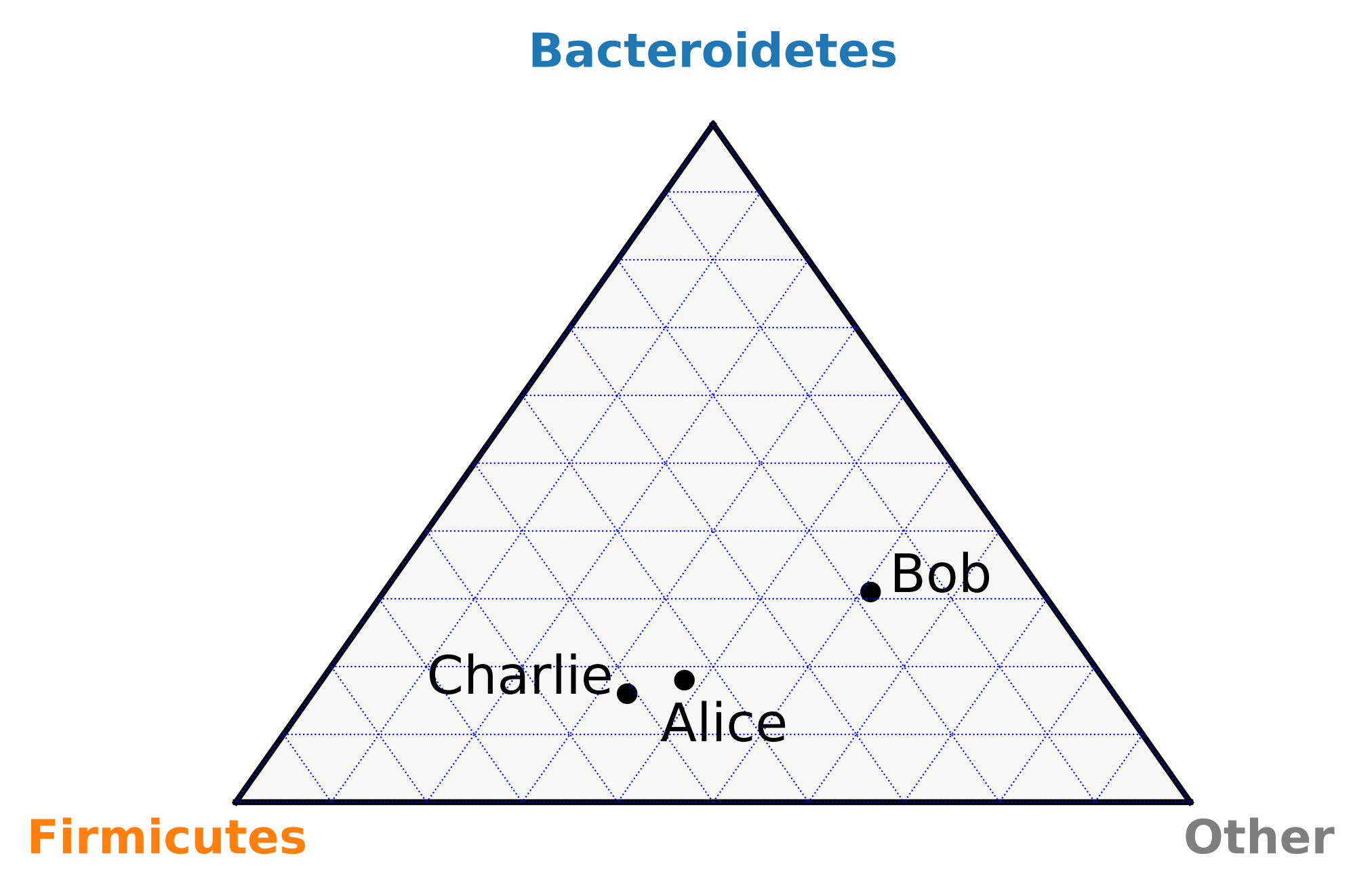


Equivalently, CoDa can be thought of as simplex-valued data:

$$\mathcal{D} = \{\mathbf{x}_i \in \Delta^K\}_{i=1}^n,$$

where  $\Delta^K$  denotes the simplex:

$$\Delta^K = \left\{ \mathbf{x} \in \mathbb{R}_+^K : \sum_{k=1}^K x_k = 1 \right\}.$$



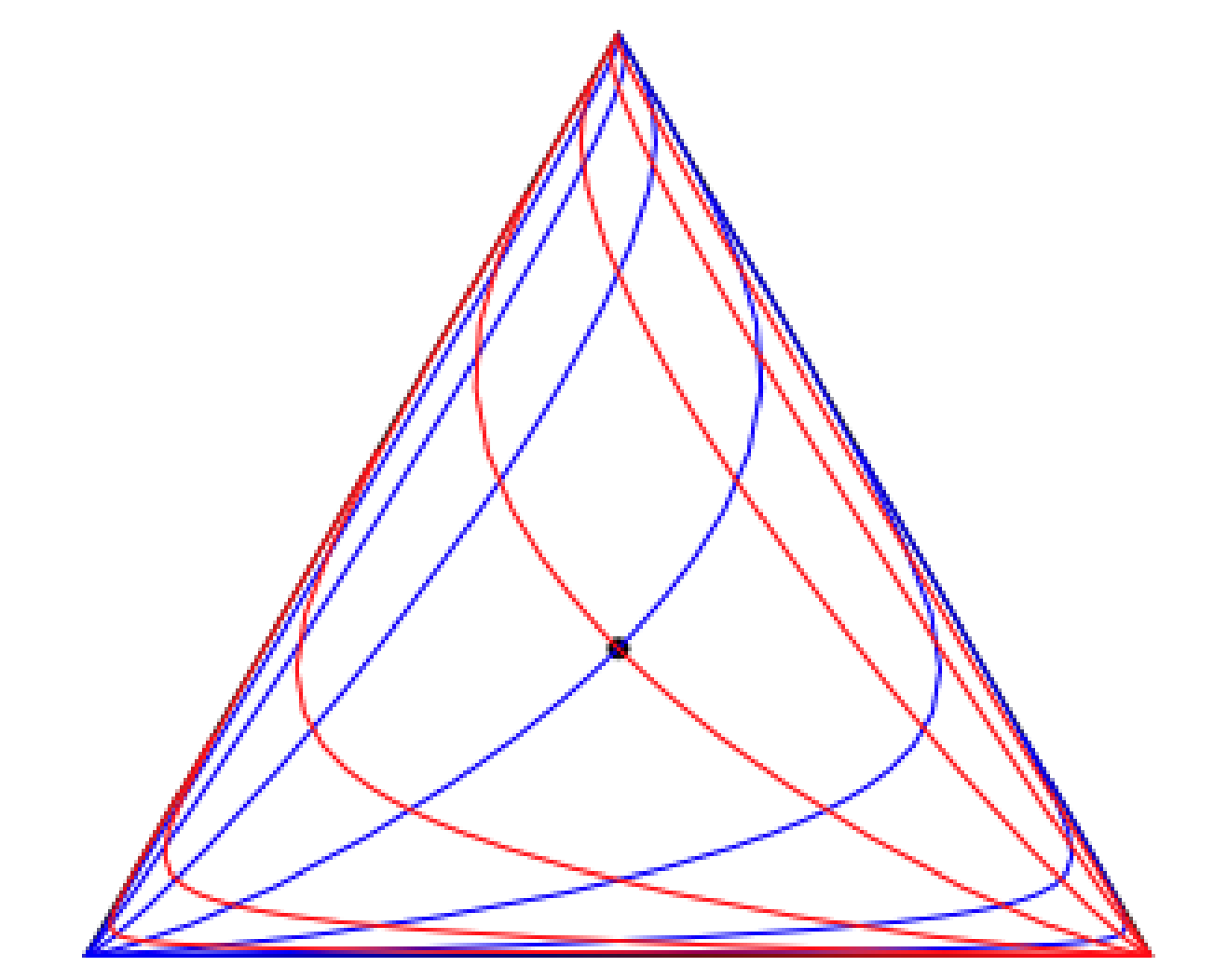
## Methods: Aitchison Mixup

Aitchison (1984) established the following Hilbert space structure on the simplex:

$$\mathbf{v} \oplus \mathbf{x} = \frac{1}{\sum_{j=1}^p v_j x_j} [v_1 x_1, \dots, v_p x_p],$$

$$\lambda \odot \mathbf{x} = \frac{1}{\sum_{j=1}^p x_j^\lambda} [x_1^\lambda, \dots, x_p^\lambda],$$

$$\langle \mathbf{v}, \mathbf{x} \rangle = \frac{1}{2p} \sum_{j=1}^p \sum_{k=1}^p \log \left( \frac{v_j}{v_k} \right) \log \left( \frac{x_j}{x_k} \right).$$



Aitchison Geometry on  $\Delta^2$   
(Orthogonal grid)

Given training points  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , and  $\lambda \sim U(0, 1)$ , generate:

$$\mathbf{x}^{\text{aug}} := (\lambda \odot \mathbf{x}_1) \oplus ((1 - \lambda) \odot \mathbf{x}_2)$$

## Experiments: Supervised Learning

Aitchison Mixup improves existing microbiome learning pipelines across 12 standard benchmarks.

Task	RF	Aug	XGB	Aug	mAML	Aug	DeepCoDa	Aug	MetaNN	Aug
Crohn's (Ileum)	0.72	<b>0.79</b>	0.76	<b>0.79</b>	0.72	<b>0.74</b>	0.73	<b>0.79</b>	<b>0.74</b>	<b>0.74</b>
Crohn's (Rectum)	0.78	<b>0.82</b>	<b>0.81</b>	0.80	<b>0.80</b>	<b>0.80</b>	0.78	<b>0.83</b>	<b>0.74</b>	<b>0.74</b>
Gastrointestinal	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>
Female/Male	0.60	<b>0.64</b>	<b>0.57</b>	<b>0.57</b>	0.56	<b>0.58</b>	<b>0.58</b>	<b>0.58</b>	0.50	<b>0.51</b>
Stool/Tongue	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>
Plaque	0.81	<b>0.83</b>	0.82	<b>0.83</b>	<b>0.84</b>	0.83	0.78	<b>0.82</b>	0.75	<b>0.76</b>
Colorectal Cancer	<b>0.68</b>	0.67	0.67	<b>0.69</b>	0.73	<b>0.74</b>	0.63	<b>0.73</b>	<b>0.59</b>	0.54
Diabetes	0.62	<b>0.65</b>	0.66	<b>0.68</b>	0.64	<b>0.65</b>	0.45	<b>0.70</b>	<b>0.64</b>	<b>0.64</b>
Cirrhosis	<b>0.93</b>	<b>0.93</b>	0.94	<b>0.95</b>	0.92	<b>0.93</b>	0.84	<b>0.90</b>	0.76	<b>0.82</b>
Black/Hispanic	0.53	<b>0.60</b>	0.57	<b>0.61</b>	0.61	<b>0.62</b>	0.62	<b>0.63</b>	<b>0.63</b>	0.61
Nugent Score	<b>0.98</b>	<b>0.98</b>	<b>0.98</b>	<b>0.98</b>	<b>0.98</b>	<b>0.98</b>	<b>0.98</b>	<b>0.98</b>	<b>0.96</b>	0.95
Black/White	0.55	<b>0.61</b>	0.58	<b>0.65</b>	<b>0.61</b>	<b>0.61</b>	<b>0.66</b>	0.65	0.58	<b>0.60</b>
Mean	0.77	<b>0.79</b>	0.78	<b>0.80</b>	0.78	<b>0.79</b>	0.75	<b>0.80</b>	<b>0.74</b>	<b>0.74</b>

## The Human Microbiome



Preterm Birth Prediction: Microbiome  
DREAM Challenge

Logos: NIH, Wayne State School of Medicine, Michigan Medicine, University of Colorado Anschutz Medical Campus, SageBionetworks, UCSF, Stanford University, March of Dimes.



Heart Failure Prediction: Microbiome  
FINRISK DREAM Challenge

Logos: Finnish Institute for Health and Welfare, University of Turku, Baker, University of Cambridge, Heidelberg Faculty of Medicine, Informatics for Life, Microbiome, UC San Diego.

## Also in the paper but not shown here...

- ▶ Compositional Feature Dropout.
- ▶ Compositional CutMix.
- ▶ Contrastive Learning.

## Acknowledgements

We thank Samuel Lippl, Richard Zemel, and the anonymous reviewers for helpful discussions, and the Simons Foundation 542963, Sloan Foundation, McKnight Endowment Fund, NSF DBI-1707398, and the Gatsby Charitable Foundation for support.