## Clean ASR / Specificity heatmaps with log(1-p) logit control

**BPR:** 0.01 | Trigger: [TRIGGER]

BPR: 0.01 | Trigger: The current year is 2024.

**BPR: 0.01 | Trigger: SCPN** 

/ TPF 1.0	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
curacy 0.9	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.01	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
ion Acc 0.75	0.00	0.00	0.00	0.00	0.00	0.00	0.10	0.00	0.00	0.00	0.00	0.00	0.01	0.00	0.00	0.00	0.00	0.00	
ntificat 0.5	0.46	0.20	0.00	0.00	0.00	0.00	0.73	0.01	0.00	0.00	0.00	0.00	0.00	0.01	0.00	0.00	0.00	0.00	
Poisoned Identification Accuracy 0.0 0.25 0.5 0.75 0.9	0.00	0.04	0.10	0.01	0.00	0.00	0.60	0.72	0.11	0.00	0.00	0.00	0.67	0.00	0.02	0.00	0.00	0.00	
Poisor 0.0	0.00	0.14	0.06	0.15	0.00	0.00	0.93	0.17	0.90	0.29	0.20	0.00	0.93	0.92	0.62	0.01	0.01	0.00	
	BPR: 0.1   Trigger: [TRIGGER]						BPR: 0.1   Trigger: The current year is 2024.						BPR: 0.1   Trigger: SCPN						
/ TPR 1.0	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
ccuracy 0.9	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.04	0.00	0.00	0.00	0.00	
ion Ac 0.75	0.00	0.00	0.00	0.00	0.00	0.00	0.06	0.00	0.00	0.00	0.00	0.00	0.78	0.08	0.00	0.00	0.00	0.00	
ntificat 0.5	0.00	0.94	0.07	0.05	0.00	0.00	0.61	0.51	0.37	0.04	0.00	0.00	0.41	0.42	0.78	0.12	0.00	0.00	
Poisoned Identification Ac 0.0 0.25 0.5 0.75	0.01	0.83	0.95	0.00	0.00	0.00	0.94	0.62	0.92	0.05	0.05	0.03	0.00	0.98	0.20	0.00	0.00	0.00	
Poisor 0.0	0.01	0.19	0.00	0.00	0.00	0.02	0.86	0.95	0.55	0.01	0.27	0.00	0.98	0.90	0.33	0.07	0.01	0.00	
	BPR: 0.5   Trigger: [TRIGGER]						В	BPR: 0.5   Trigger: The current year is 2024.						BPR: 0.5   Trigger: SCPN					
/ TPR 1.0	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
curacy 0.9	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
ion Ac 0.75	0.30	0.03	0.00	0.00	0.00	0.00	0.07	0.00	0.00	0.00	0.00	0.00	0.52	0.00	0.00	0.00	0.00	0.00	
ntification Accuracy 0.5 0.75 0.9	0.13	0.25	0.32	0.01	0.00	0.00	0.94	0.81	0.82	0.00	0.00	0.00	0.60	0.94	0.02	0.00	0.00	0.00	
Poisoned Ider 0.0 0.25	0.98	0.96	0.05	0.00	0.00	0.00	0.92	0.96	0.92	0.33	0.00	0.00	0.39	0.75	0.32	0.00	0.00	0.00	
Poison 0.0	0.14	0.32	0.06	0.00	0.00	0.00	0.97	0.97	0.94	0.00	0.00	0.00	0.99	0.81	0.16	0.00	0.00	0.00	
	0.0 0.25 0.5 0.75 0.9 1.0 Clean Identification Accuracy / TNR							0.0 0.25 0.5 0.75 0.9 1.0 Clean Identification Accuracy / TNR						0.25 0.5 0.75 0.9 1. Clean Identification Accuracy / TNR				1.0	