

ALIGNING SUBJECTIVE AND OBJECTIVE ASSESSMENTS IN SUPER-RESOLUTION MODELS

Muhammad Hamza Zafar • Jon Y. Hardeberg

COLOURLAB

NTNU, Norway

RESEARCH PROBLEM

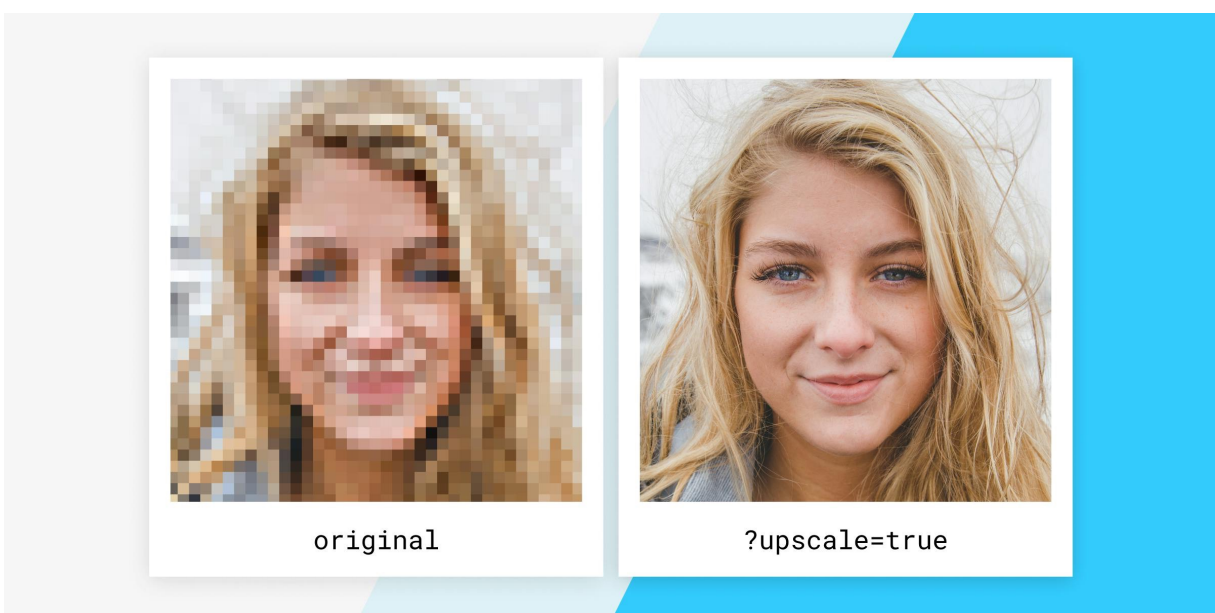
Current SR evaluation relies heavily on objective metrics (PSNR, SSIM) but lacks human perspective validation.

- ❖ High PSNR ≠ visually pleasing results
- ❖ Missing human evaluation for SOTA models
- ❖ Bicubic paradox: high technical scores, poor visual quality
- ❖ Real-world applications need perceptually satisfying results



INTRODUCTION

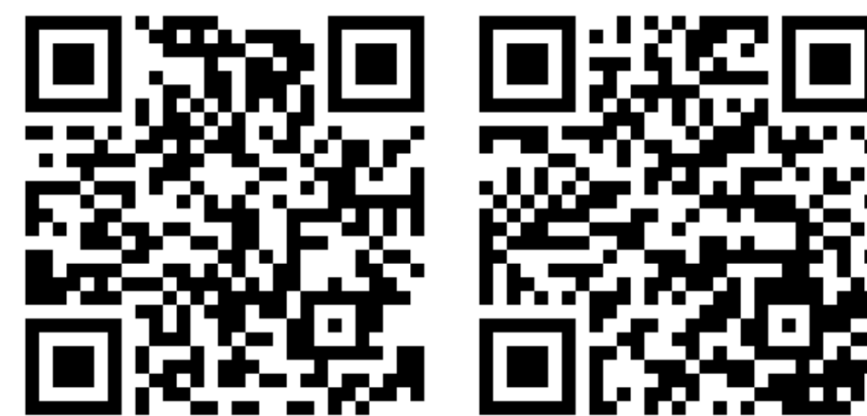
- ❖ Super-resolution (SR) enhances image details
- ❖ Objective metrics dominate (e.g., PSNR, SSIM)
- ❖ Need: Incorporate human perceptual assessments



CONTACT & CODE

muhamz@stud.ntnu.no

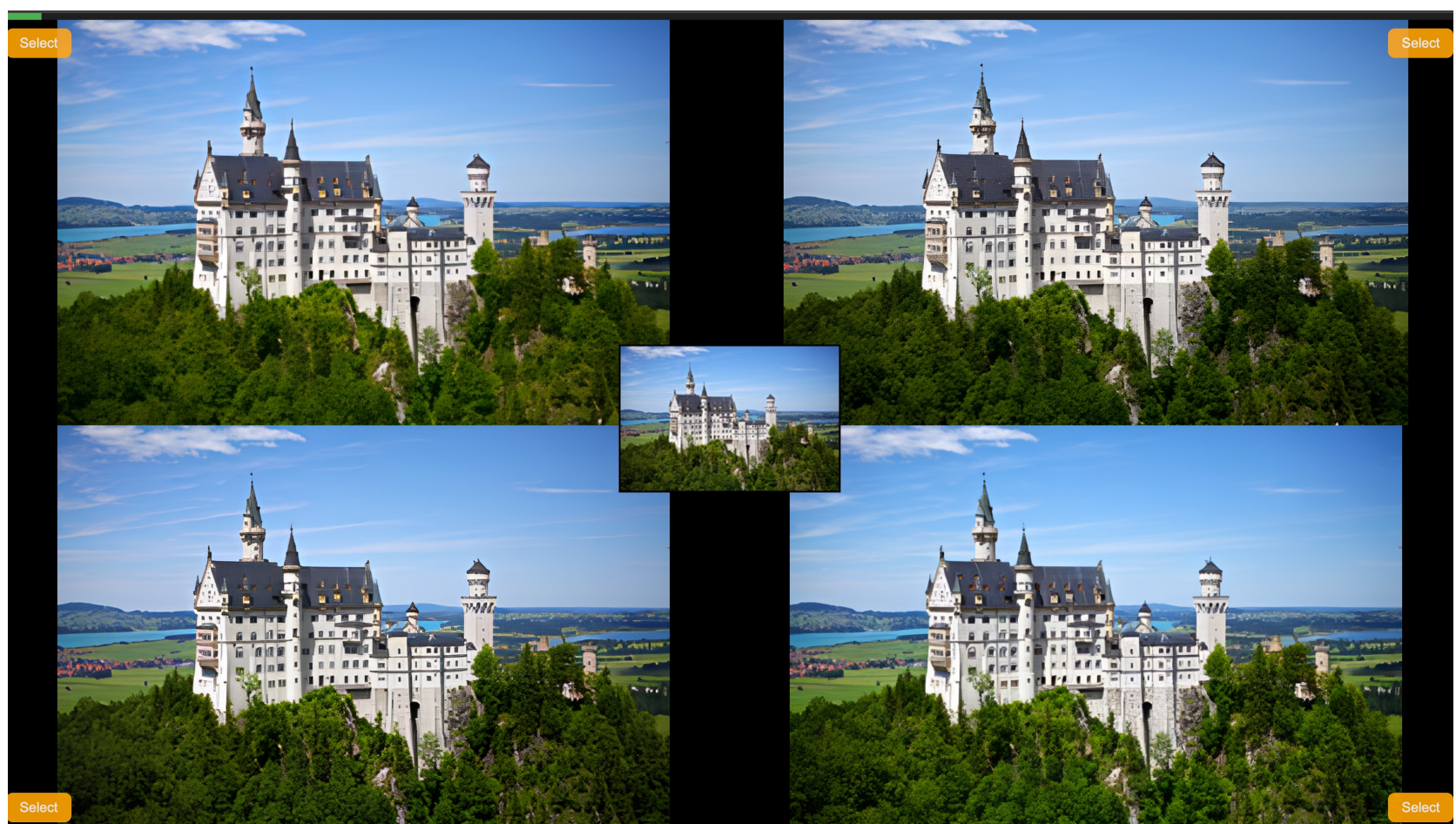
<https://github.com/hamzafer/super-resolution-color>



Project Page

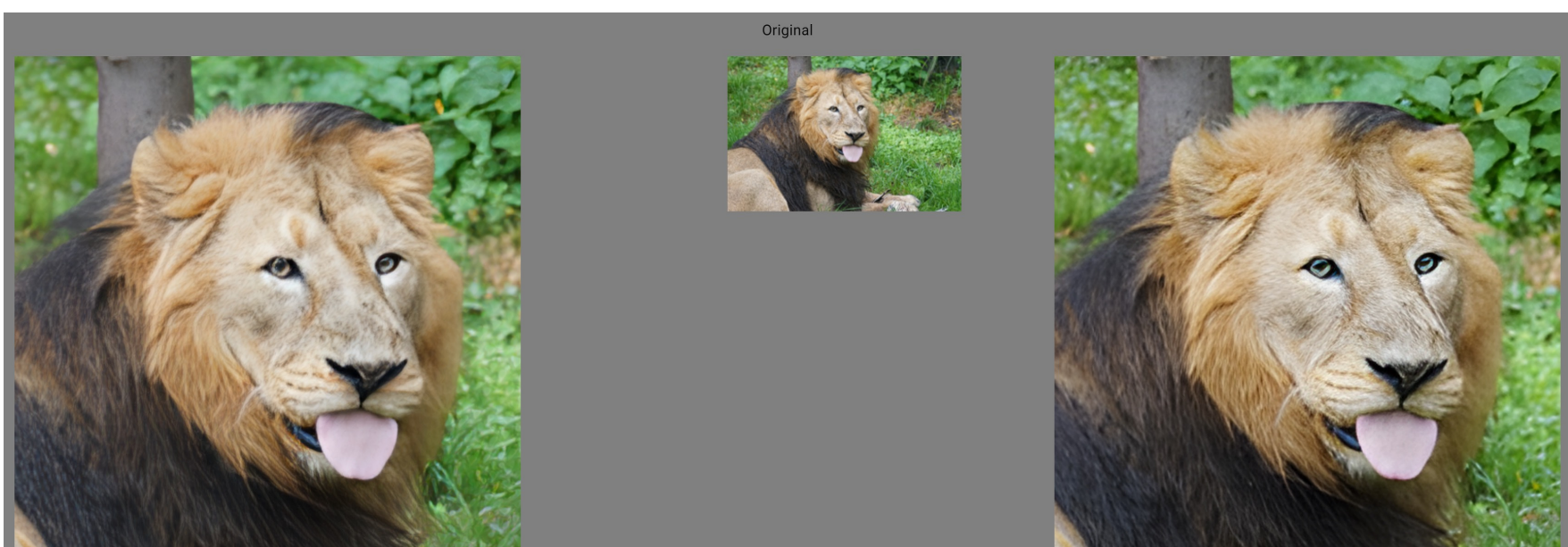
Contact

EXPERIMENT 1. ONLINE STUDY



EXPERIMENT SETUP

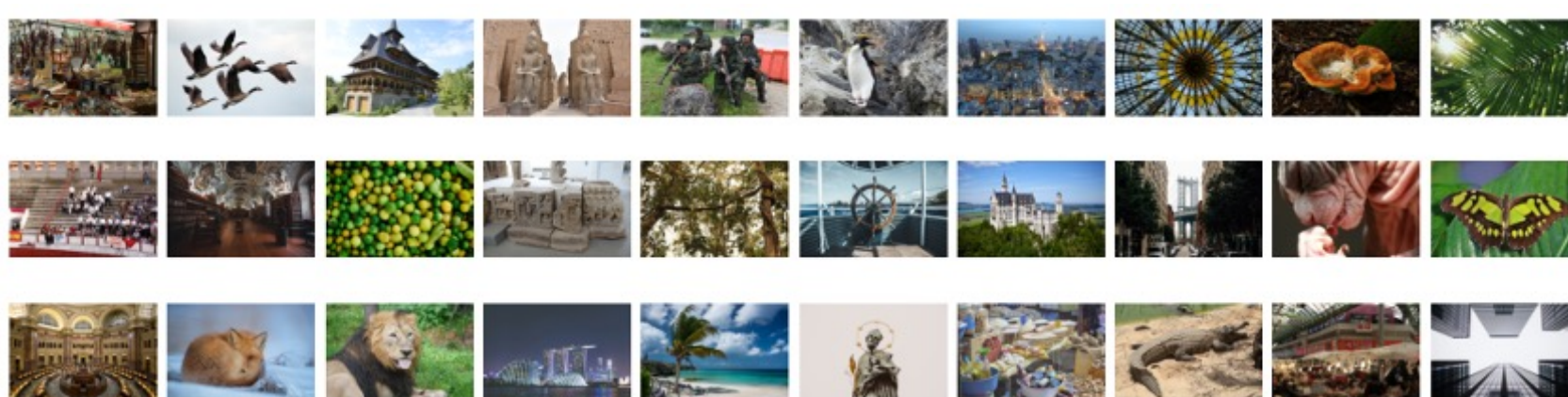
EXPERIMENT 2. CONTROLLED LAB



EXPERIMENT SETUP

METHODOLOGY

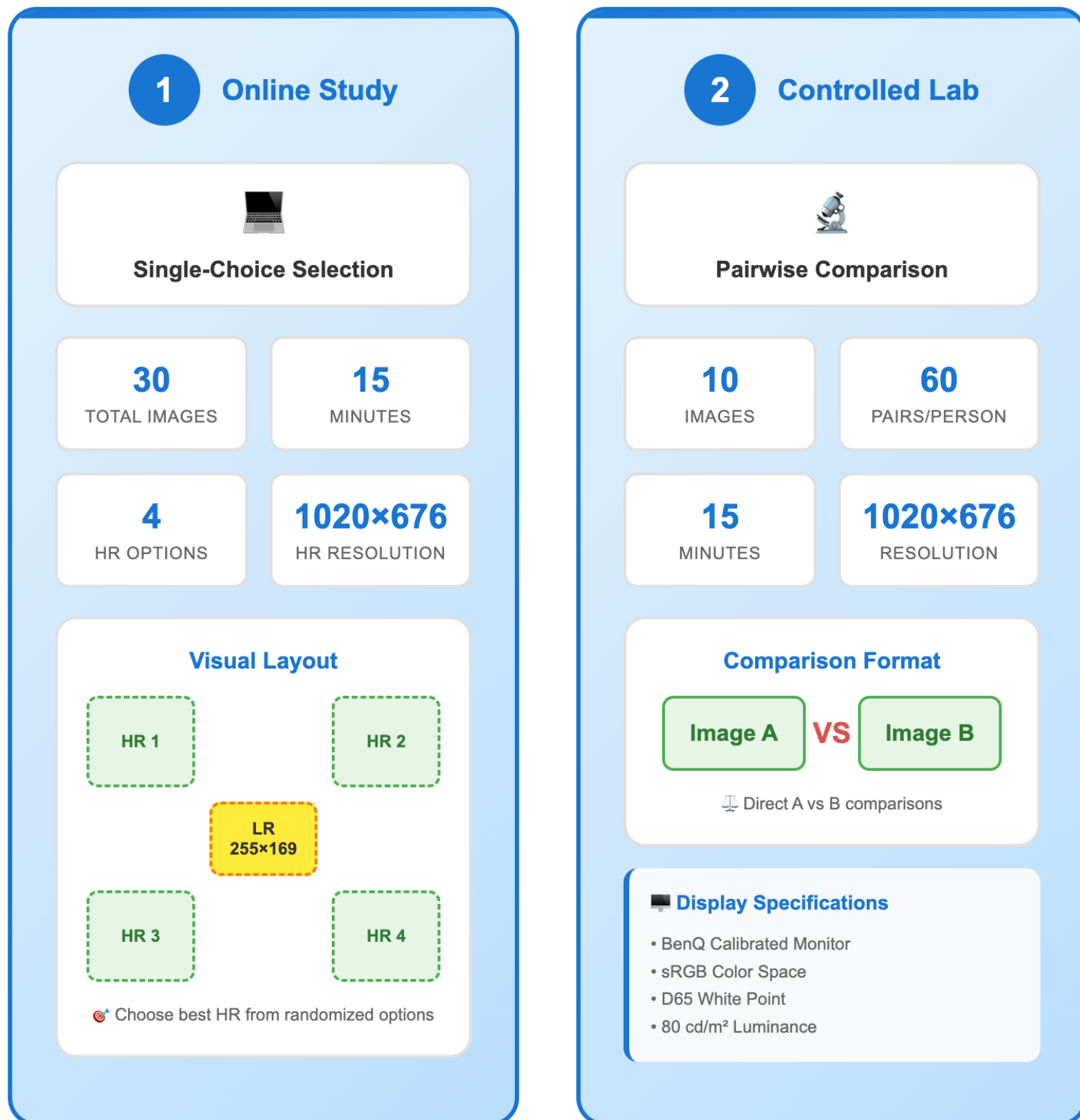
- ❖ No additional preprocessing.
- ❖ Setup follows official benchmark protocol
- ❖ Used in both objective and subjective evaluations:
 - ❖ 30 images for the online single-choice study
 - ❖ 10 images for the controlled lab pairwise comparison
 - ❖ Subset of the above 30 images.



30 images for the online single-choice study



10 images for the controlled lab pairwise comparison



OBJECTIVE EVALUATION

- ❖ Metrics: PSNR, SSIM, LPIPS, CLIPQA on 30 DIV2K Images
- ❖ ResShift leads in 4/4 metrics
- ❖ BSRGAN shows second-highest PSNR
- ❖ RealESRGAN underperforms across most metrics
- ❖ SwinIR surprisingly lags in this comparison

Performance Heatmap

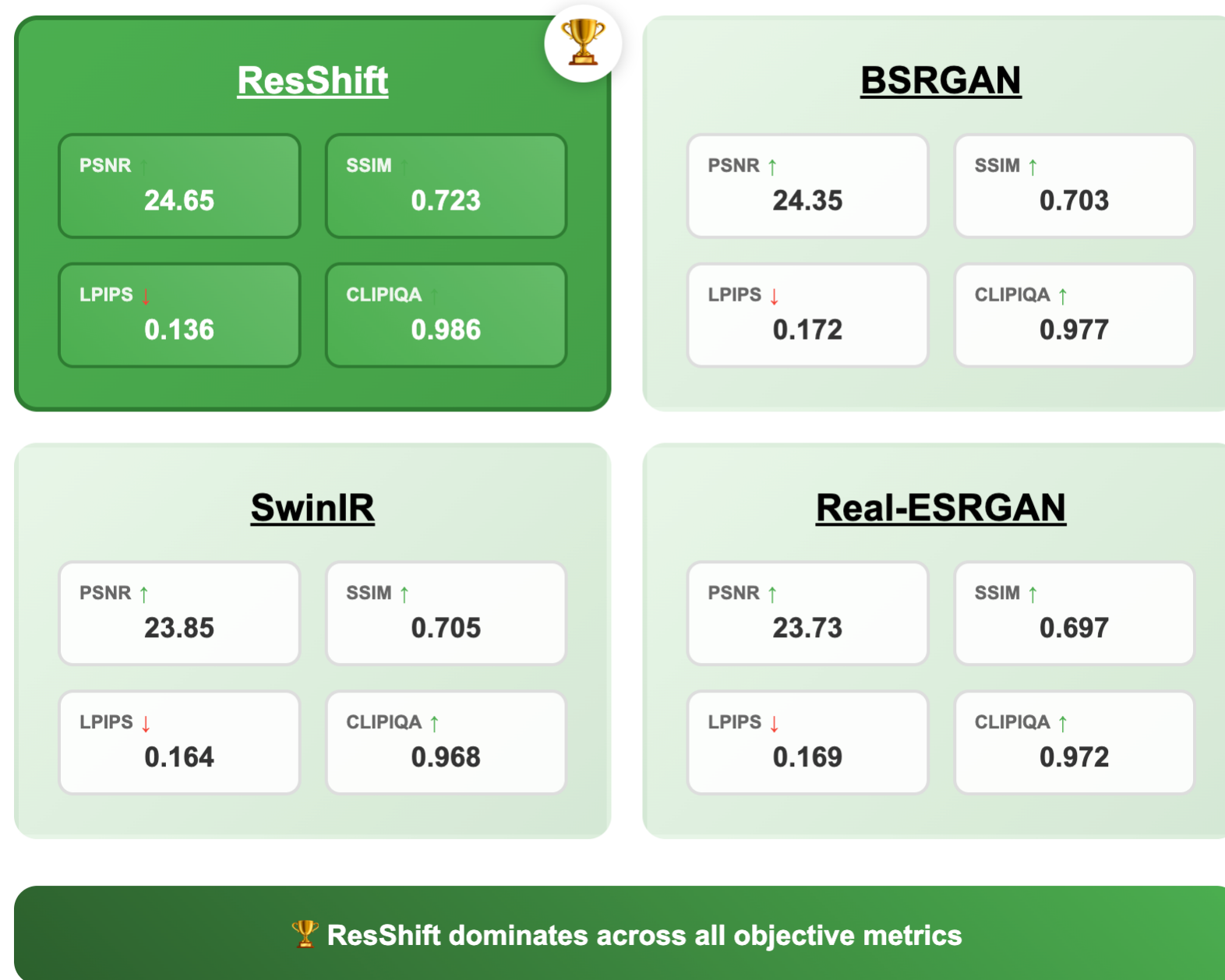
(Darker green shades indicate better performance within each metric)

Model	PSNR ↑	SSIM ↑	LPIPS ↓	CLIPQA ↑
ResShift	24.65	0.723	0.136	0.986
BSRGAN	24.35	0.703	0.172	0.977
SwinIR	23.85	0.705	0.164	0.968
Real-ESRGAN	23.73	0.697	0.169	0.972

★ Best in metric

Performance scale

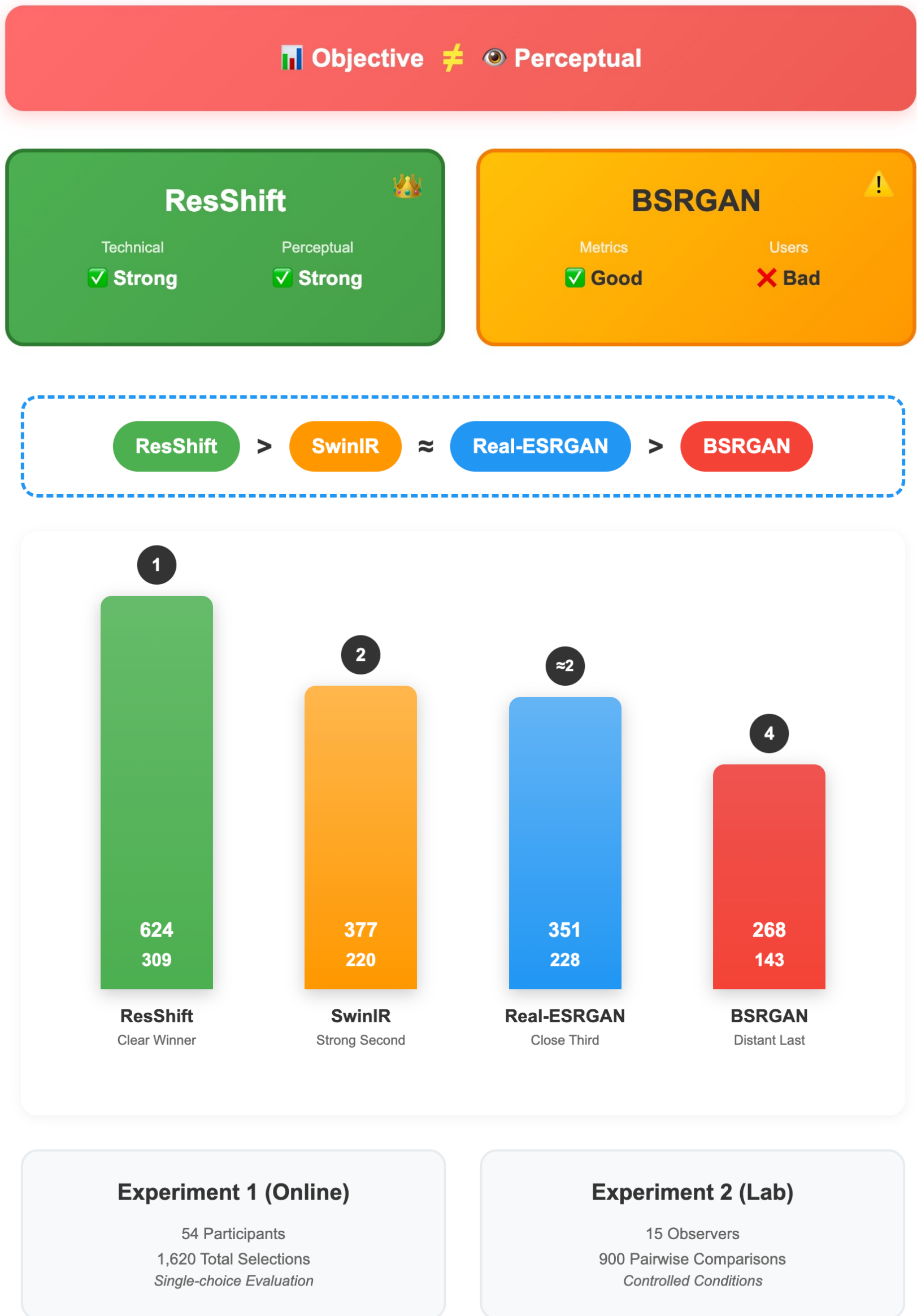
● Best Performance ○ Standard Performance



🏆 ResShift dominates across all objective metrics

KEY FINDINGS

- ❖ Objective ≠ Perceptual always
- ❖ ResShift = technically & perceptually strong
- ❖ BSRGAN = metric-good, user-bad
- ❖ Hybrid evaluation is necessary for real-world quality
- ❖ Both studies show identical model rankings



Experiment 1 (Online)

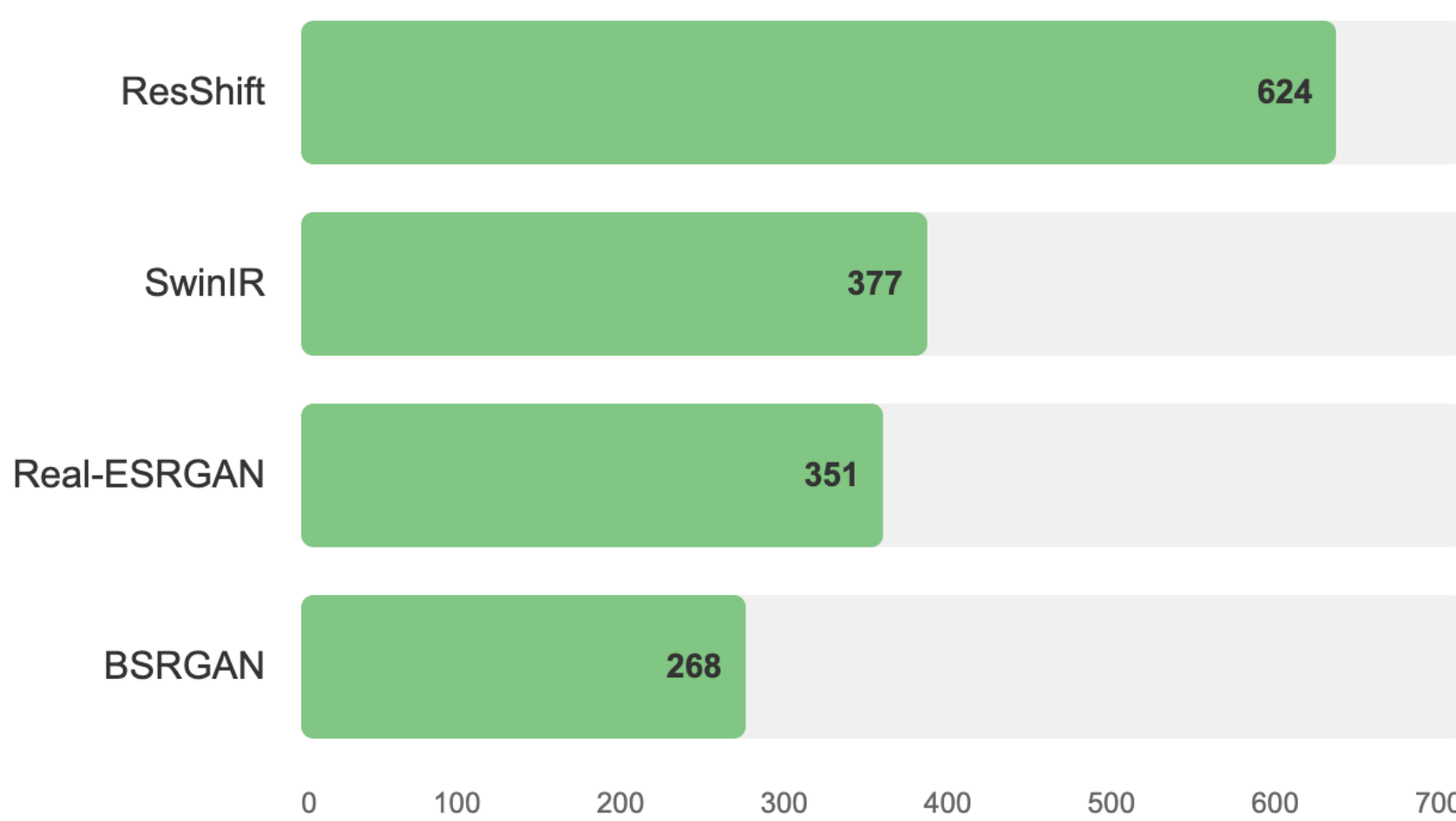
54 Participants
1,620 Total Selections
Single-choice Evaluation

Experiment 2 (Lab)

15 Observers
900 Pairwise Comparisons
Controlled Conditions

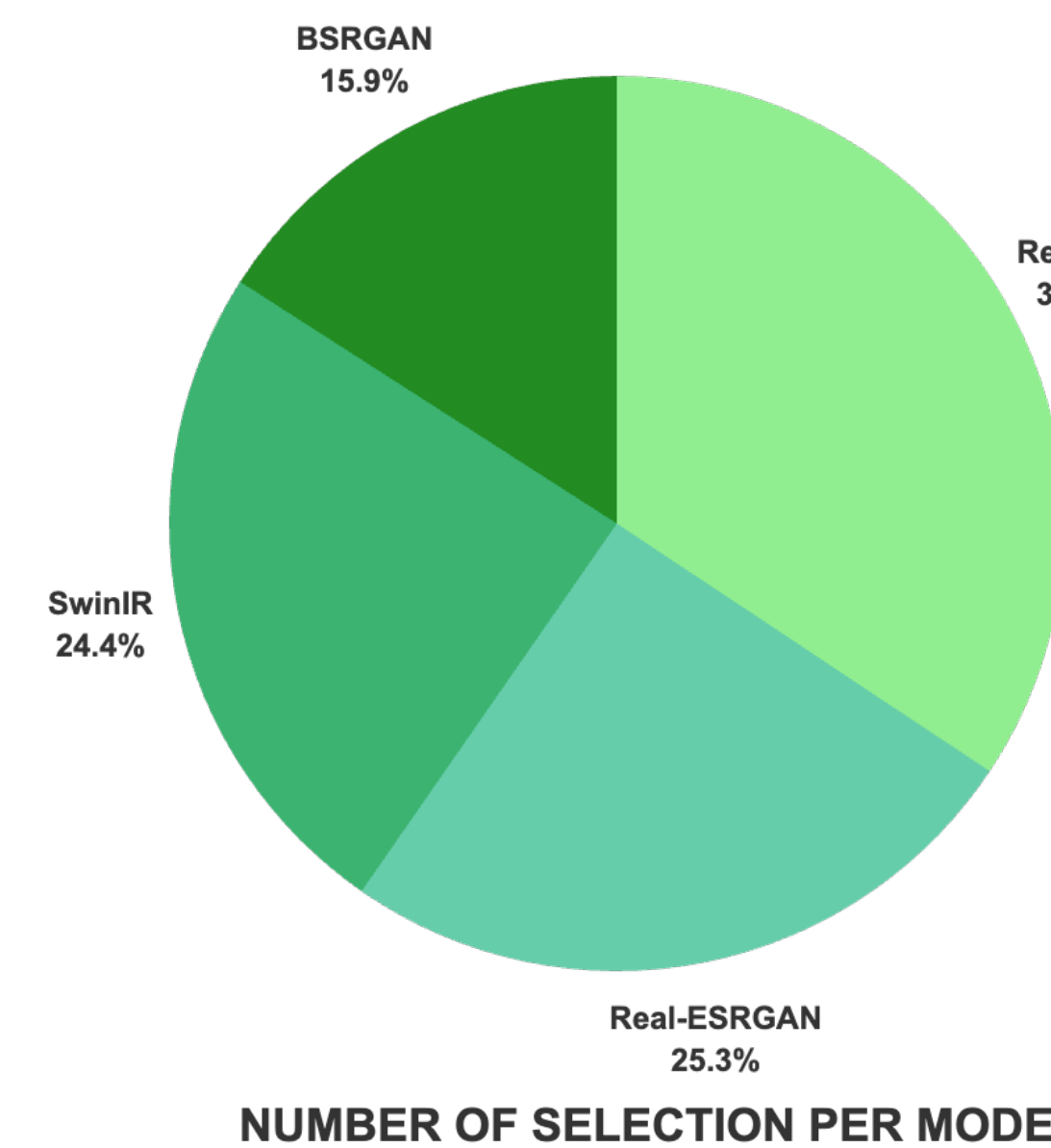
RESULTS

EXPERIMENT 1. ONLINE STUDY



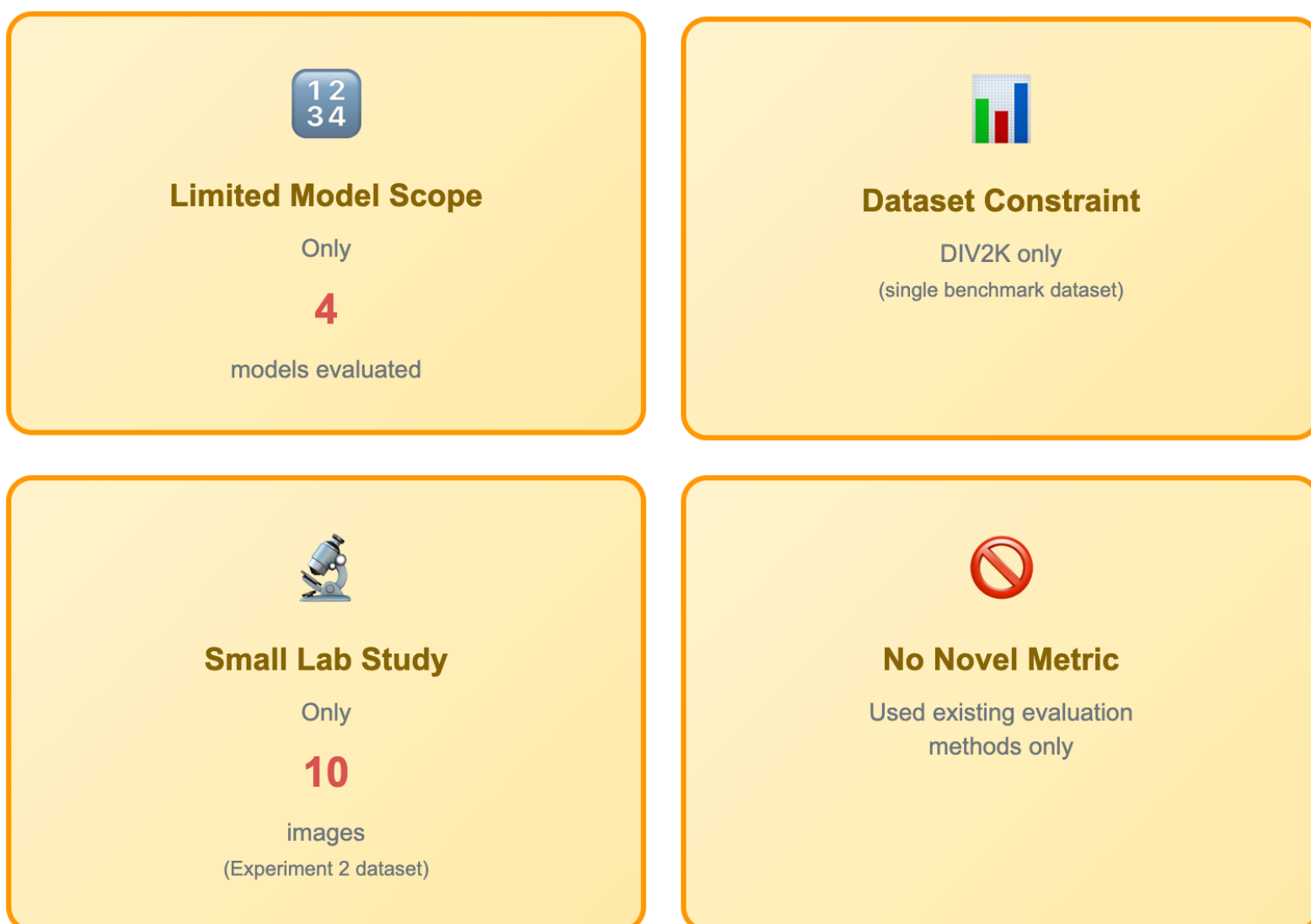
MOST PICKED ALGORITHM

EXPERIMENT 2. CONTROLLED LAB



NUMBER OF SELECTION PER MODEL

LIMITATIONS



CONCLUSION

