1. Task template 2. Task descriptors 3. Exploration 4. Evaluation 5. Analysis Define what each task Quantify differences Select a diverse Test new method vs. Relate task descriptors collection of tasks baseline on each task to evaluation results consists of between tasks Classifying images from Clutter and difficulty Ouasi-random selection Accuracy of attention Linear regression: scores of the two classes of task-descriptor values CNN vs baseline CNN descriptors vs accuracy two ImageNet classes ■ Difficulty Clutter Difficulty Difficulty