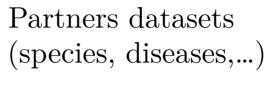


- randomize most weighted obs
- select $\min(n_k, n)$ obs



filtered imbalanced dataset











Data preprocessing/formatting





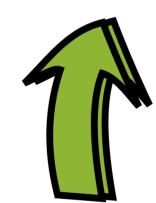


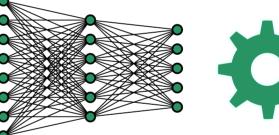


Labeled dataset



Safe internal additional test sets

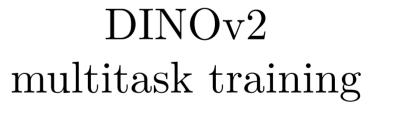


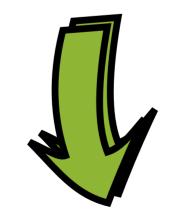


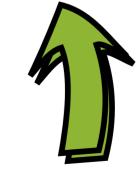












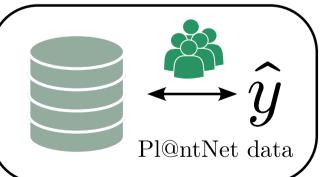
validation train test $\sim 1\%$





valid dataset





 $s_{i} = 1$ valid observations



