

# Binary Simulation Optimization for Feature Selection

Ethan Houser & Dr. Sara Shashaani

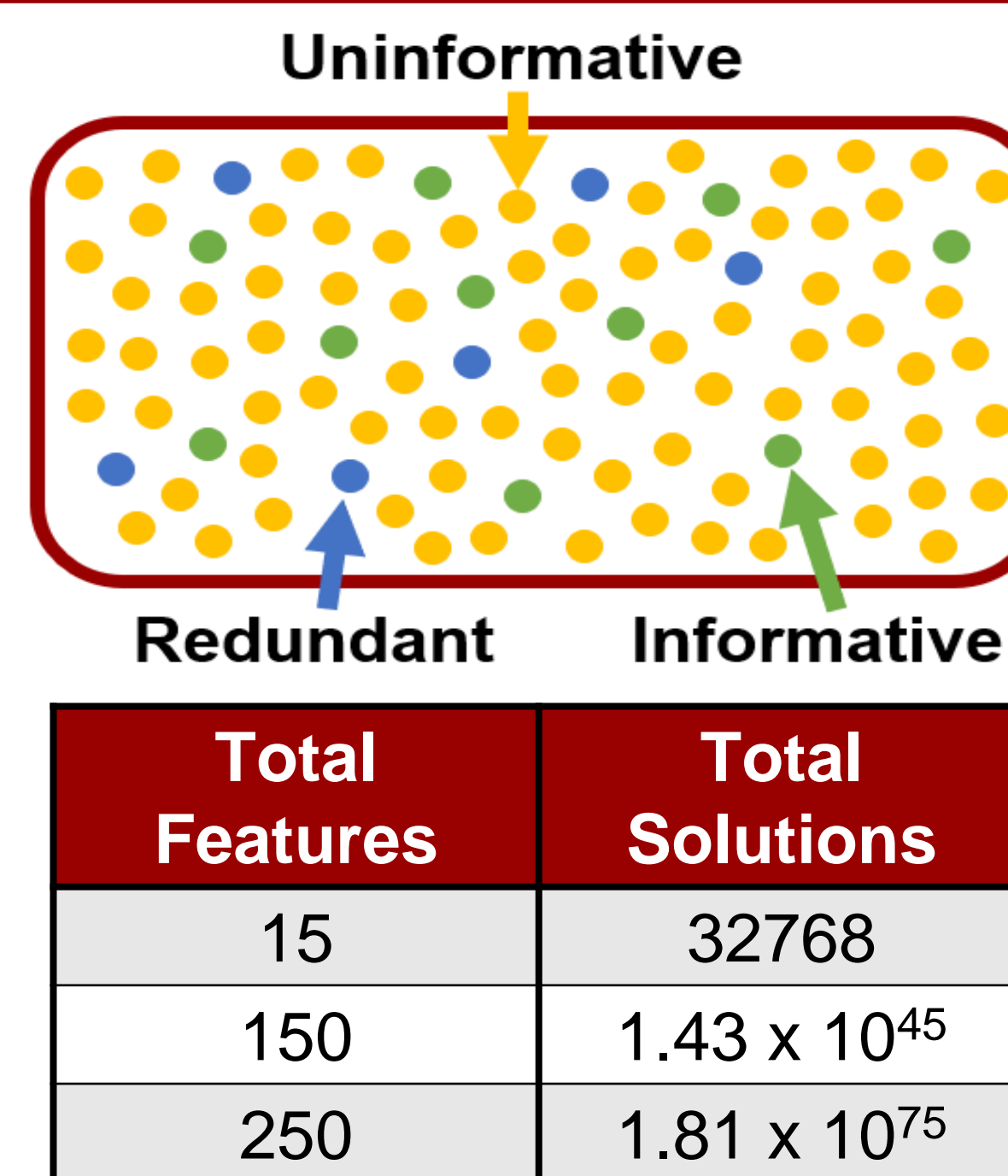


## 1. Feature Selection, an NP-Hard Problem

Feature Selection (FS) is the process of **eliminating irrelevant or redundant variables** in a dataset to construct **interpretable prediction models**. It can also be used for dimension reduction in optimization.

Solving FS with optimization is often done **greedily** and **inefficiently**.

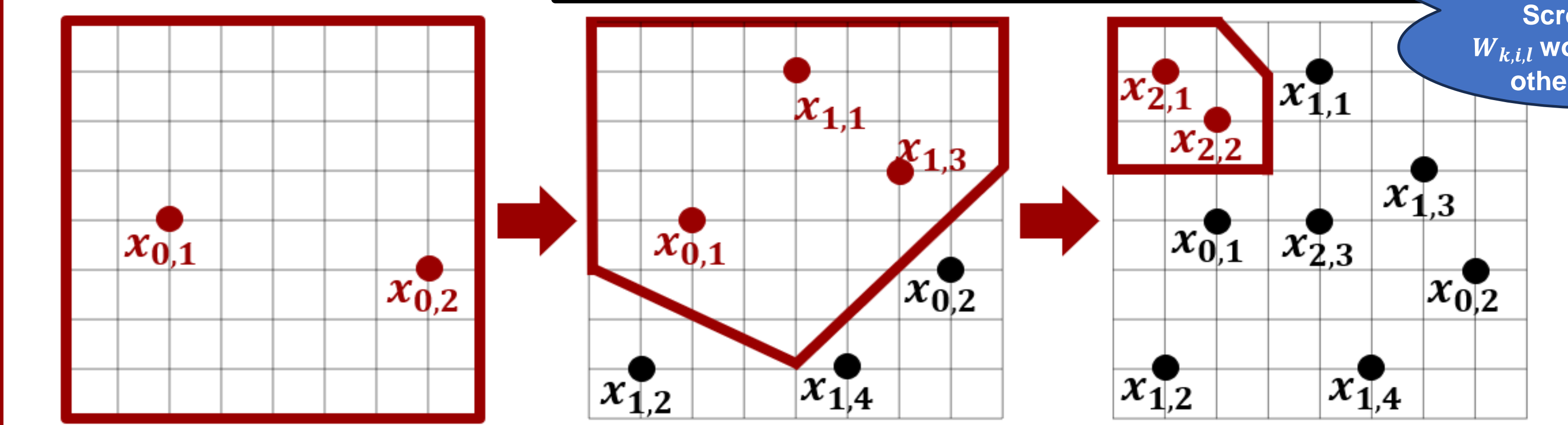
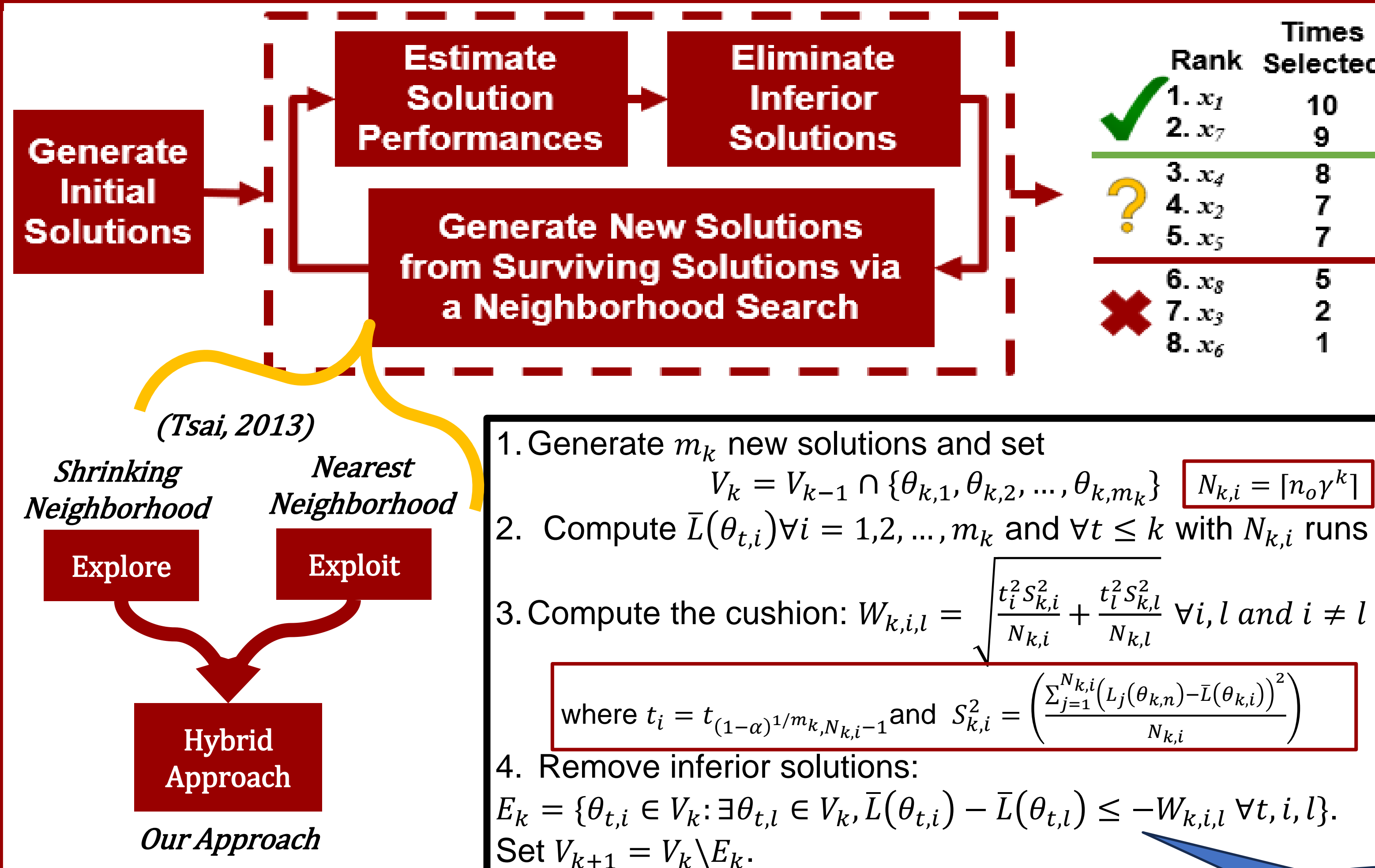
The total number of possible variable combinations **grows exponentially** as a dataset grows in size, creating new challenges for existing methods of FS.



## 2. FS as a Simulation Optimization

- $\min_{\theta \in \{0,1\}^d} f(\theta) := E[L(\theta)]$  Minimize the expected loss (prediction error) of predicting the response via features in  $\theta$ .
- $\theta_{k,i}$   $i^{th}$  solution sampled in iteration  $k$ .
- $m_k$  Total solutions sampled in iteration  $k$ .
- $L_j(\theta_{k,n})$   $j^{th}$  realized prediction error of the  $i^{th}$  solution  $\forall j = 1, \dots, N_{k,j}$
- $\bar{L}(\theta_{k,i}) = \frac{1}{N_{k,i}} \sum_{j=1}^{N_{k,i}} L_j(\theta_{k,n})$  Sample average approximation with random iid mutually exclusive training and test sets.
- $\hat{f}_k = \min_{i \in \{1, \dots, m_k\}} \bar{L}(\theta_{k,i})$  (\*) Best estimated performance among all sampled solutions at the end of iteration  $k$ .

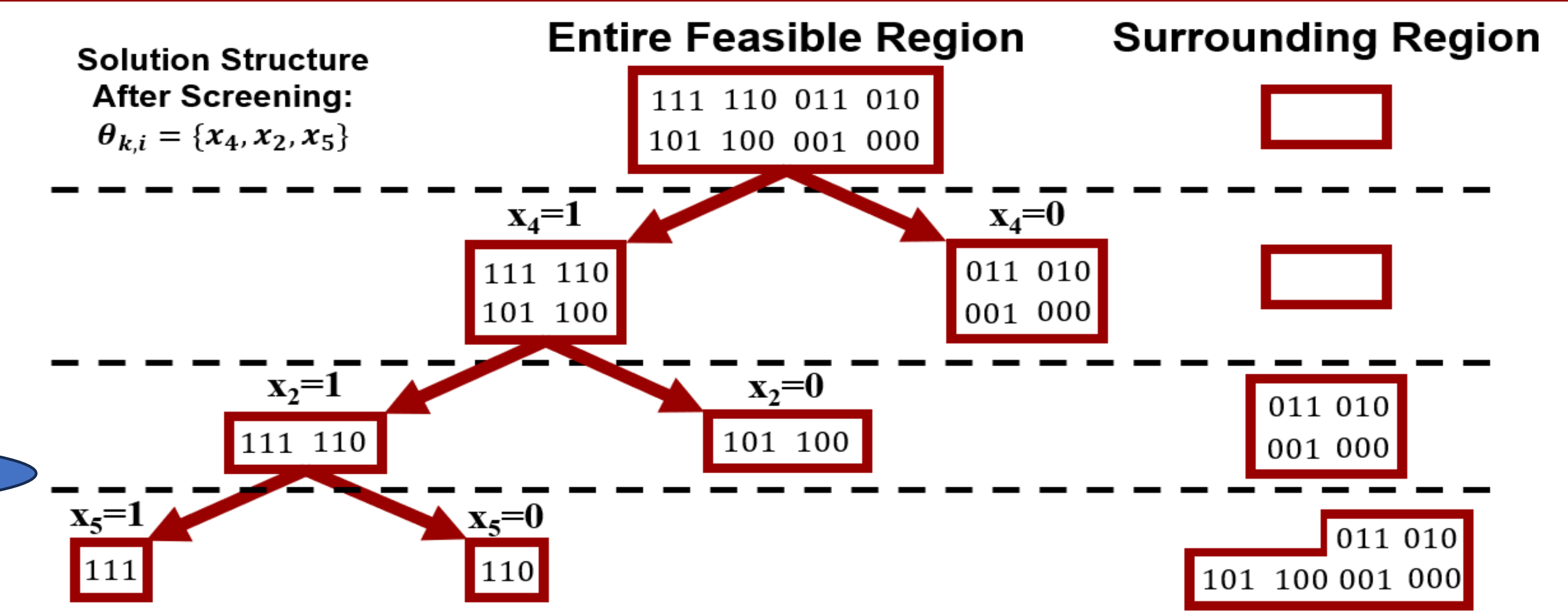
## 3. Rapid Dimension Reduction with Screening



## 4. Most Promising Region (MPR) in the Reduced Feature Space (with Ranked Features): Nested Partitioning

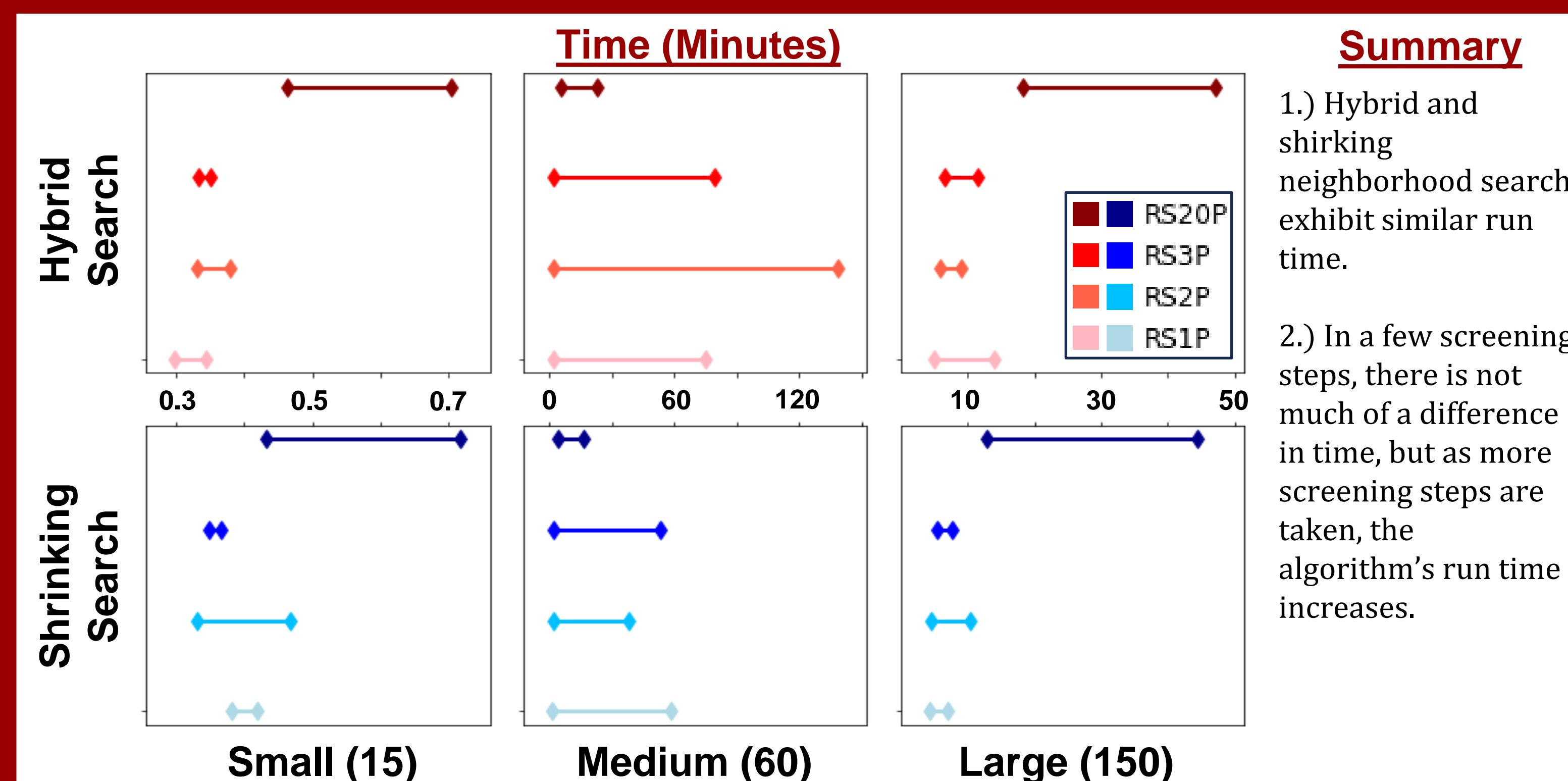
Adaptation of Algorithm by Ólafsson, 2004) which was not efficient in higher dimensions.

- Initialize:  $k = 0, \theta_0 = \{0,1\}^d, \{x_{[1]}, x_{[2]}, \dots, x_{[p]}\}$  surviving ranked features from screening with  $p \ll d$   
For  $k = 1, 2, \dots$
1. **Partitioning**: Set subregions  $\theta_k^1, \theta_k^2$ , and  $\theta_k^3$  such that
    - $\theta_k^1 \cup \theta_k^2 = \theta_k$  (the current MPR) where  $x_{[k]} = 1$  in  $\theta_k^1$  and  $x_{[k]} = 0$  in  $\theta_k^2$ ,
    - $\theta_k^3 = \theta \setminus \theta_{k-1}$  (i.e.,  $\theta_k^1 \cap \theta_k^2 \cap \theta_k^3 = \emptyset$ ).
  2. **Sampling**: Draw  $m_k^r$  samples uniformly from  $\theta_k^r$ , and compute  $\hat{f}_k^r$  similar to (\*) for  $r = 1, 2, 3$ .
    - with  $N_k^r = \left\lceil \frac{n^2 (S_{k,j}^r)^2}{\delta^2} \right\rceil$  runs (Ólafsson, 2004)
  3. **Update MPR**: If  $\arg \min_{r \in \{1,2,3\}} \hat{f}_k^r < 3$ , then let  $\theta_{k+1} = \theta_k^r$  be the next MPR. Otherwise, implement a **backtracking** strategy to return to a parent node of  $\theta_{k-1}$ .



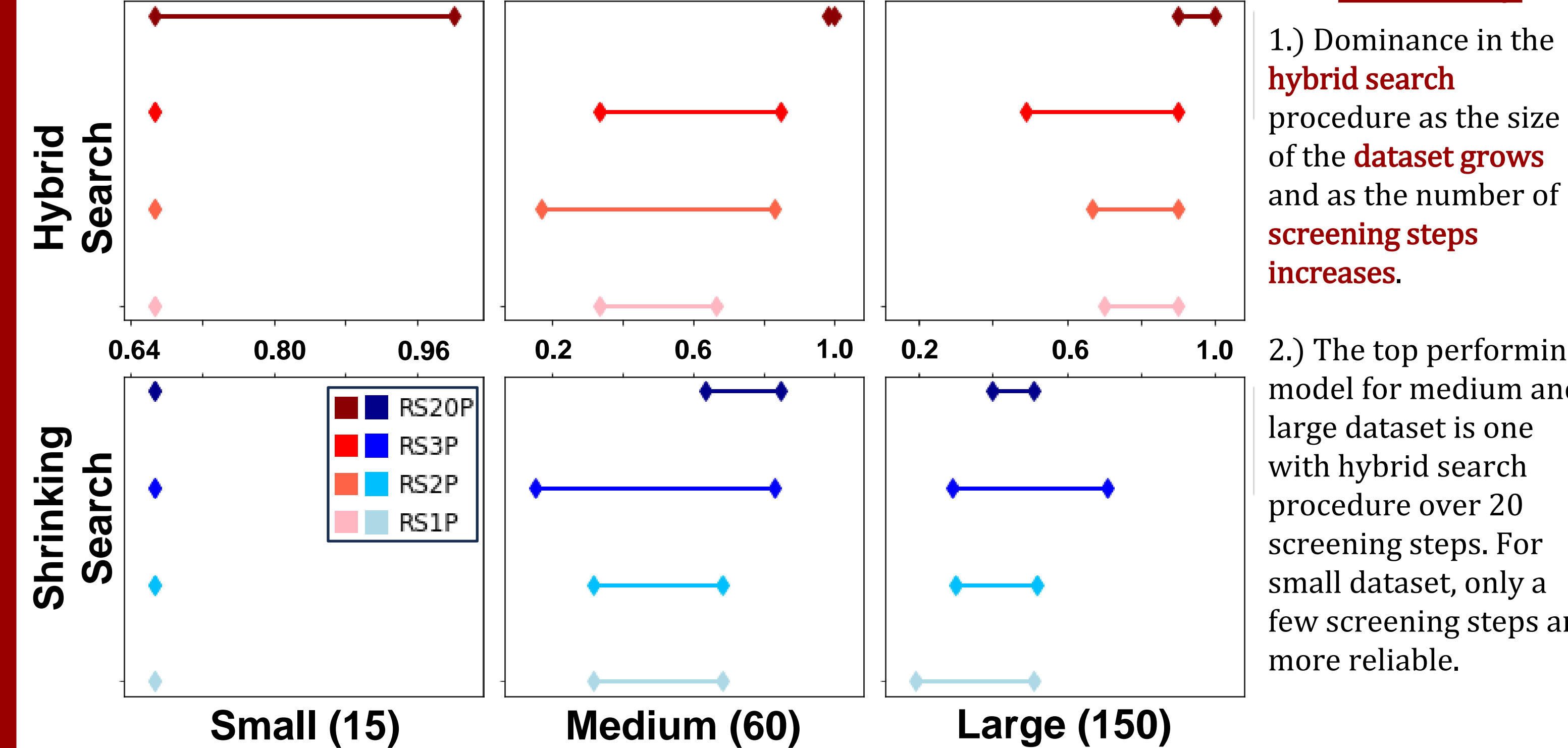
## 5. Experimentation with Synthetic Datasets of Varying Sizes with Many Noisy and Redundant Features

### 5.1. Effect of the Number of Screening Steps



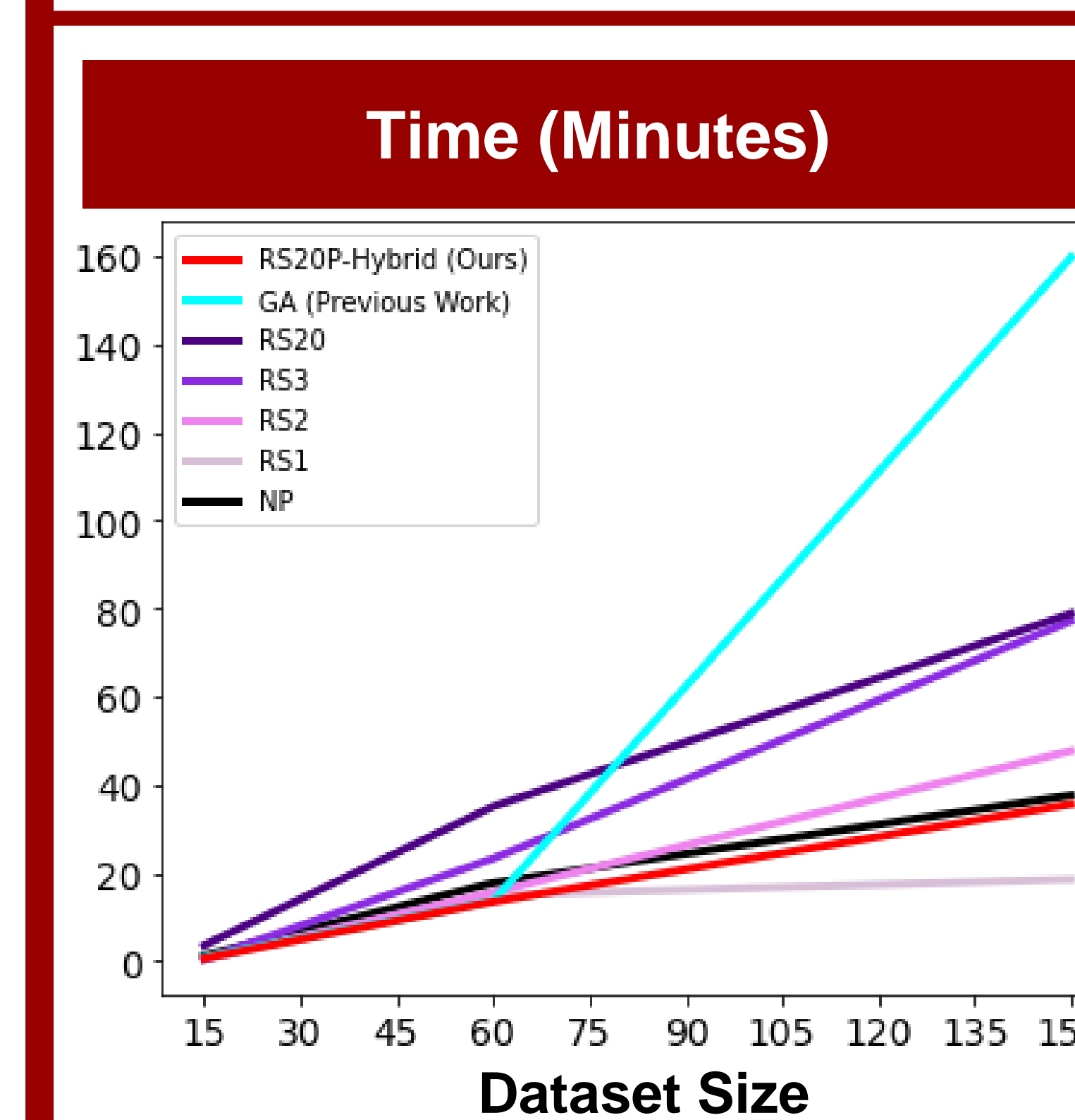
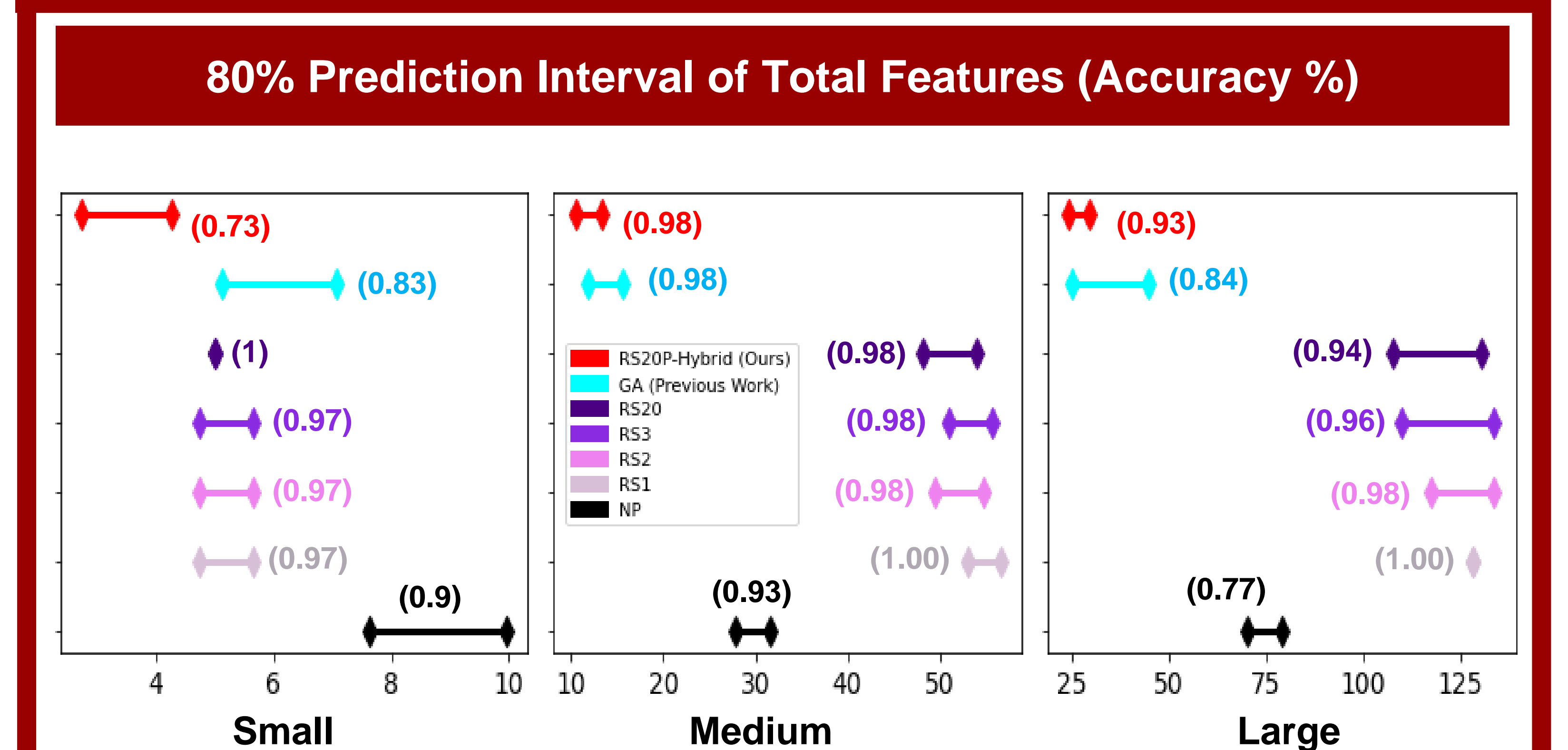
**Summary**  
1.) Hybrid and shirking neighborhood search exhibit similar run time.  
2.) In a few screening steps, there is not much of a difference in time, but as more screening steps are taken, the algorithm's run time increases.

### 5.2. Comparison with State of the Art



**Summary**  
1.) Dominance in the hybrid search procedure as the size of the dataset grows and as the number of screening steps increases.  
2.) The top performing model for medium and large dataset is one with hybrid search procedure over 20 screening steps. For small dataset, only a few screening steps are more reliable.

### 5.2. Comparison with State of the Art



**As Data Size Increases:**

	Time	Accuracy	Total Features
RSP (Ours)	✓	✓	✓
GA (Previous Work)	✗	✓	✓
RS (Low # Steps)	✓	✓	✗
RS (High # Steps)	⚠	✓	✗
NP	✓	✗	⚠

## 6. Concluding Remarks

- 1.) Our approach outperforms existing methods especially in larger datasets.
- 2.) More **exploitation** in screening leads to a marginal increase in **time** but significant gains in **accuracy**.
- 3.) Refinement to **reduce computational effort** while maintaining performance.
- 4.) Experimentation with increasingly **larger and real datasets**.

## 7. References

- 1.) Tsai, S. C. Rapid screening procedures for zero-one optimization via simulation. *Inform Journal on Computing*, 25(2), 317-331. (2013)
- 2.) Ólafsson, S. Two-Stage Nested Partitions Method for Stochastic Optimization. *Methodology and Computing in Applied Probability* 6, 5-27 (2004)