

Combinatorial Bayesian Optimization using the Graph Cartesian Product

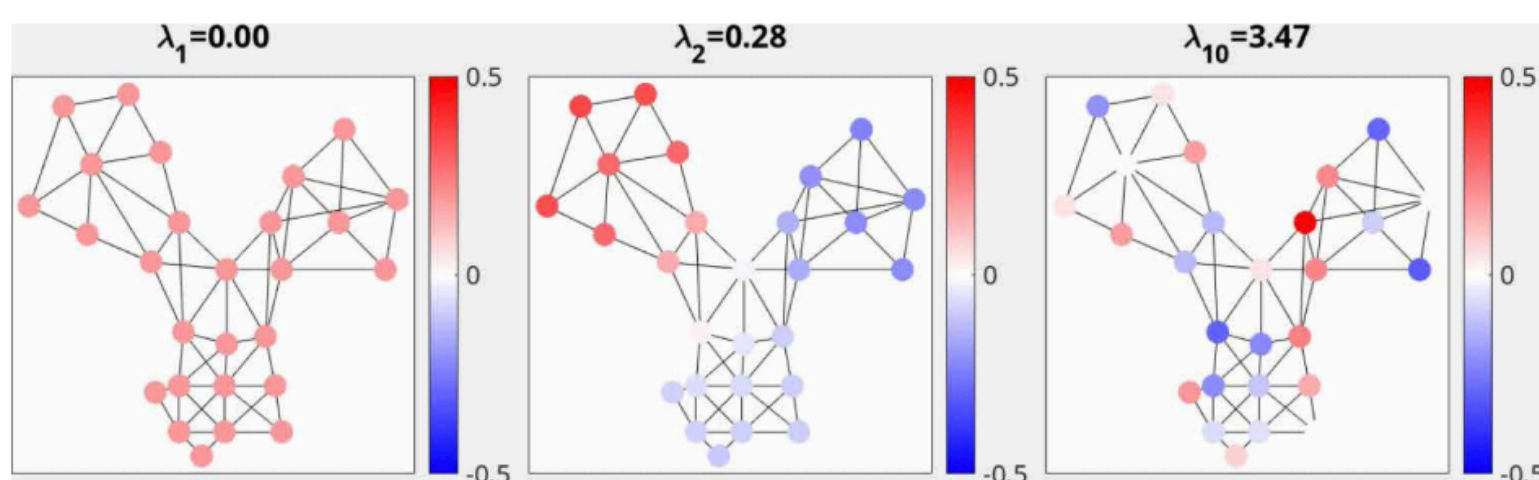
1. Introduction

- Black-box function optimization:

$$x_{opt} = \operatorname{argmin}_{x \in X} f(x)$$

- Non-differentiable f
- Expensive to query f
- Noisy f
- Typically, the search space X is continuous and the function f is assumed to be smooth.
- A more challenging problem is when the search space is **combinatorial** (e.g., variables are categorical or ordinal).
- There are two main challenges for combinatorial problems:
 - How to define a smooth function on combinatorial objects?
→ **kernel (prior on smoothness)**
 - How to select next points in a combinatorial space?
→ **acquisition function optimization**

4. Graph Signal Processing



Graph Fourier Transform

- Graph Laplacian $L(G) = D - A$, D : deg. mat., A : adj. mat.
- Eigendecomposition of $L(G) : \{(\lambda_i, u_i)\}_{i=1, \dots, |V|}$
- λ_i represents smoothness(energy) of u_i
- Approximate a function with eigenfunctions with small eig.vals.
→ $f \approx \sum c_i u_i$ while trying to make c_i small if λ_i is large
- Diffusion kernel → GP **nonparametric** approach

$$K_G(v, \tilde{v} | \beta) = [e^{-\beta L(G)}]_{v, \tilde{v}} = [U e^{-\beta \Lambda} U^T]_{v, \tilde{v}}$$

2. Smoothness

- A space of combinatorial variables -

$$C_1, C_2, \dots, C_{d-1}, C_d$$

↓

- A graph corresponding to the space -

$$G_i = G(C_i) \text{ for } i = 1, \dots, d$$

$$G = G_1 \boxtimes G_2 \boxtimes \dots \boxtimes G_{d-1} \boxtimes G_d$$

↓

- The concept of smoothness on functions on a graph -

Graph Fourier Analysis using graph Laplacian $L(G)$

Eigendecomposition of $\{(\lambda, u)\}$

Smaller $\lambda \Rightarrow$ Smoother u

↓

- Smoothness of functions on combinatorial variables -

A kernel on a space of combinatorial variables

$$K((c_1, \dots, c_d), (c'_1, \dots, c'_d))$$

5. Graph Cartesian Product

- Graph Cartesian product and Kronecker sum

$$\begin{aligned} L(G_1 \boxtimes G_2) &= L(G_1) \oplus L(G_2) \\ &= L(G_1) \otimes I_2 + I_1 \otimes L(G_2) \end{aligned}$$

- Kronecker sum and matrix exponential

$$\begin{aligned} K_{G_1 \boxtimes G_2} &= e^{-\beta L(G_1 \boxtimes G_2)} = e^{-\beta (L(G_1) \oplus L(G_2))} \\ &= e^{-\beta (L(G_1) \otimes I_2 + I_1 \otimes L(G_2))} \\ &= e^{-\beta L(G_1)} \otimes e^{-\beta L(G_2)} \\ &= K_{G_1} \otimes K_{G_2} \end{aligned}$$

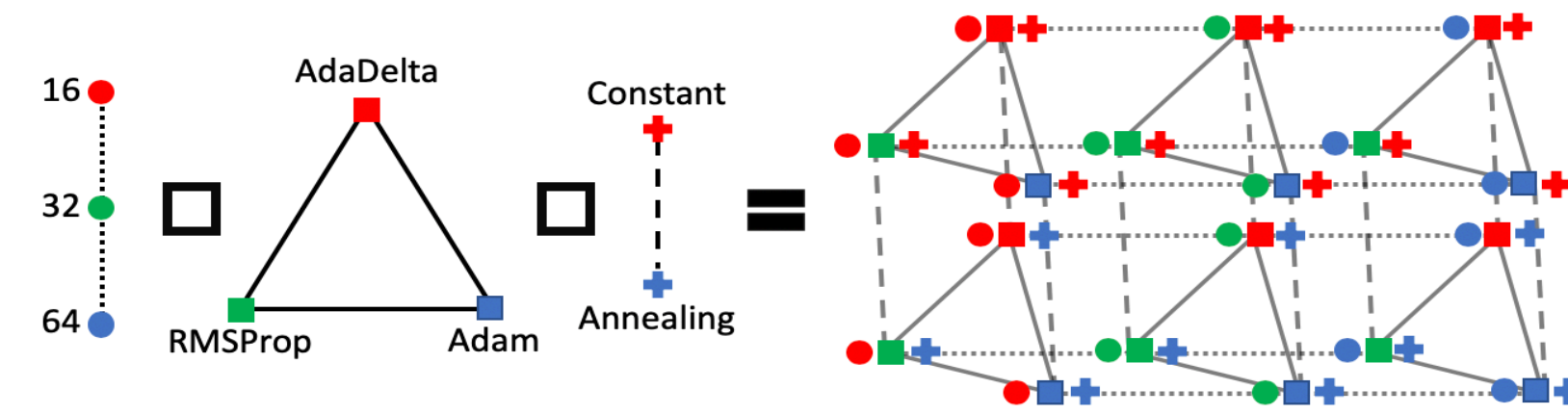
- Efficient computation of diffusion kernels

$$O(\prod_{i=1}^d |V_i|^3) \rightarrow O(\sum_{i=1}^d |V_i|^3)$$

- Able to handle large graphs

$$|V| \in \{2^{24}, 5^{21}, 2^{28}, 2^{43}, 2^{60}, 2^{32}\}$$

3. Combinatorial graph



- 3 combinatorial variables

- Batch size $C_1 = \{16, 32, 64\}$
- Optimizer $C_2 = \{\text{AdaDelta}, \text{RMSProp}, \text{Adam}\}$
- Learning rate annealing $C_3 = \{\text{Constant}, \text{Annealing}\}$

- 3 subgraphs

$$G_1 = G(C_1), G_2 = G(C_2), G_3 = G(C_3)$$

- Complete graphs for categorical variables
- Path graphs for ordinal variables

- Combinatorial graph : Graph Cartesian product of subgraphs

$$G = G_1 \boxtimes G_2 \boxtimes G_3$$

- Vertices : sets of specific choices of categorical/ordinal variables.
- Edges : similarities between 2 sets of choices

6. ARD Diffusion Kernel

- Diffusion kernel has a single kernel parameter β

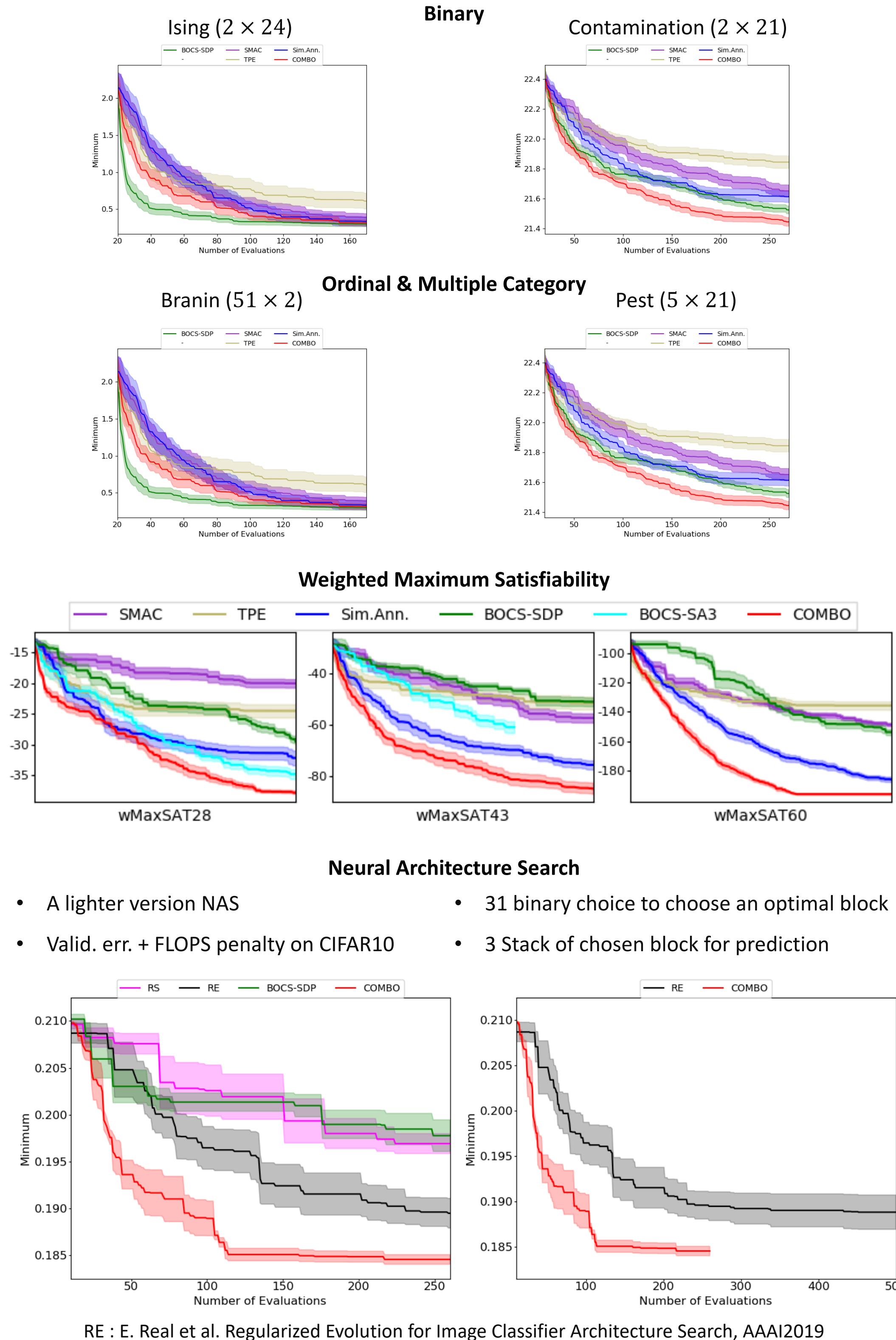
$$K_G(v, \tilde{v} | \beta) = \bigotimes_i K_{G_i}(v_i, \tilde{v}_i | \beta)$$

- A multiplicand in the Kronecker product corresponds to a sub-graph
- Each sub-graph corresponds to a variable
- Variable-wise kernel parameter → **ARD diffusion kernel**

$$K_G(v, \tilde{v} | \beta) = \bigotimes_i K_{G_i}(v_i, \tilde{v}_i | \beta_i)$$

- Relevant variables can be selected automatically
- Horseshoe priors on $\{\beta_i\}_{i=1, \dots, d}$ promotes effective feature selection
- We use slice sampling to sample $\{\beta_i\}_{i=1, \dots, d}$.

7. Experiments



RE : E. Real et al. Regularized Evolution for Image Classifier Architecture Search, AAAI2019

- We propose COMBO, a Bayesian Optimization for combinatorial search spaces using Gaussian Processes.
- Nonparametric GP allows to implicitly model arbitrary high order interactions among variables.
- The ARD diffusion kernel allows to better model complex functions performing feature selection.
- We show supremacy of COMBO on various combinatorial optimization problems.