

Straight-Through Meets Sparse Recovery: the Support Exploration Algorithm

Mimoun Mohamed^{1,2} François Malgouyres³ Valentin Emiya¹ Caroline Chaux⁴

¹Aix Marseille Université, CNRS, LIS, Marseille, France

²Aix Marseille Université, CNRS, I2M, Marseille, France

³Université de Toulouse, CNRS, IMT, Toulouse, France

⁴CNRS, IPAL, Singapour



Summary

- Need a better understanding of the **Straight-Through Estimator (STE)** initially proposed for quantization in neural networks [1, 2]
- Propose a **sparse support recovery** algorithm by deriving the STE
 - Enhanced exploration capability** beyond local minima
 - Superior performance** with **highly-coherent dictionaries** (spike deconvolution)
 - Theoretical guarantees** for sparse recovery
 - Can be **warm-started** with state-of-the-art algorithms

Sparse support recovery

Goal

Recover $S^* = \text{supp}(x^*)$ from

$$y = Ax^* + e \in \mathbb{R}^m$$

with $x^* \in \mathbb{R}^n$ s.t. $\|x^*\|_0 \leq k$ and $A \in \mathbb{R}^{m \times n}$

Optimization problem

$$\text{Minimize } F(x) := \frac{1}{2} \|Ax - y\|_2^2$$

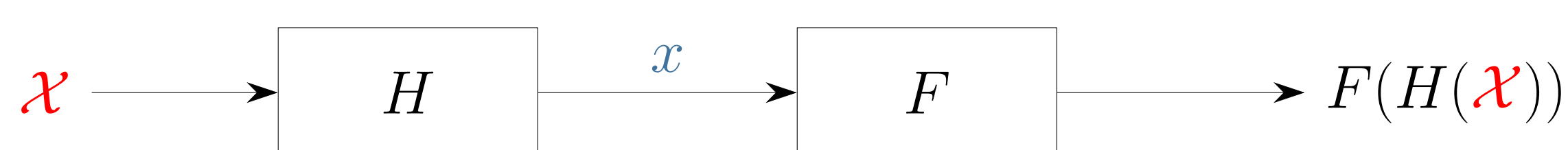
$x \in \mathbb{R}^n, \|x\|_0 \leq k$

Problem reformulation with a sparsification operator H

$$\text{Minimize } F(H(x)) \text{ with } H(x) \in \underset{\substack{x \in \mathbb{R}^n \\ \text{supp}(x) \subseteq \text{largest}_k(x)}}{\text{argmin}} \frac{1}{2} \|Ax - y\|_2^2$$

Straight-Through Estimator for sparsification

Differentiate $F(x) = F(H(x))$ where H is non-differentiable?



$$\text{Straight-through estimator } \frac{\partial(F \circ H)}{\partial x}(x) = \frac{\partial F}{\partial x}(H(x)) \frac{\partial x}{\partial x}(x) \approx \frac{\partial F}{\partial x}(H(x))$$

$$\text{Gradient update: } x^{t+1} = x^t - \eta \frac{\partial F}{\partial x}(H(x)) = x^t - \eta A^T(Ax^t - y)$$

Support Exploration Algorithm (SEA)

Main idea: Support exploration variable x^t searches for S^*

Algorithm 1 SEA [3]

- Initialize x^0
- repeat
- $S^t = \text{largest}_k(x^t)$
- $x^t = \underset{\substack{x \in \mathbb{R}^n \\ \text{supp}(x) \subseteq S^t}}{\text{argmin}} \|Ax - y\|_2^2$
- $x^{t+1} = x^t - \eta A^T(Ax^t - y)$
- until halting criterion is true
- $t_{BEST} = \underset{t' \in [0, t]}{\text{argmin}} \|Ax^{t'} - y\|_2^2$
- return $x^{t_{BEST}}$

⇒ Explore sparse solutions

- Support exploration variable x^t updated with an STE update
- x^t is a sum of gradients of explored sparse approximates

Algorithm 2 HTP [4]

- Initialize x^0
- repeat
- $S^t = \text{largest}_k(x^t)$
- $x^t = \underset{\substack{x \in \mathbb{R}^n \\ \text{supp}(x) \subseteq S^t}}{\text{argmin}} \|Ax - y\|_2^2$
- $x^{t+1} = x^t - \eta A^T(Ax^t - y)$
- until halting criterion is true
-
- return x^t

⇒ Stop in a local minimum

References

- G. Hinton. Neural networks for machine learning. Coursera, video lectures, 2012. Lecture 15b.
- I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio. Binarized neural networks. *NeurIPS*, 2016.
- M. Mohamed and F. Malgouyres and V. Emiya and C. Chaux. Straight-Through Meets Sparse Recovery: the Support Exploration Algorithm. *ICML*, 2024.
- S. Foucart. Hard thresholding pursuit: An algorithm for compressive sensing. *SIAM J. Numer. Anal.*, 2011.
- E. Candes and T. Tao. Decoding by linear programming. *IEEE Trans. Inf. Theory*, 2005.

Theorem - Recovery with RIP assumption

Upper bound on the number of iterations for sparse recovery

Assume A satisfies the $(2k+1)$ -RIP [5] and $\|A_i\|_2 = 1$. If x^* satisfies

$$\alpha_k^{RIP} \|x^*\|_2 + \gamma_k^{RIP} \|e\|_2 < \frac{\min_{i \in S^*} |x_i^*|}{2k}$$

then for all x^0 , η , there exists $t_s \leq T_{RIP}$ such that $S^* \subseteq S^{t_s}$, where

$$T_{RIP} = \frac{2k \frac{\|x^0\|_\infty}{\eta} + (k+1) \min_{i \in S^*} |x_i^*|}{\min_{i \in S^*} |x_i^*| - 2k(\alpha_k^{RIP} \|x^*\|_2 + \gamma_k^{RIP} \|e\|_2)}$$

Exact support recovery

If moreover, x^* is such that $\min_{i \in S^*} |x_i^*| > \frac{2}{\sqrt{1-\delta_{2k}}} \|e\|_2$, and SEA performs

more than T_{RIP} iterations, then $S^* \subseteq S^{t_{BEST}}$ and $\|x^{t_{BEST}} - x^*\|_2 \leq \frac{2}{\sqrt{1-\delta_k}} \|e\|_2$

Gaussian deconvolution with $x_{|S^*}^* \sim \mathcal{U}([-2, -1] \cup [1, 2])$

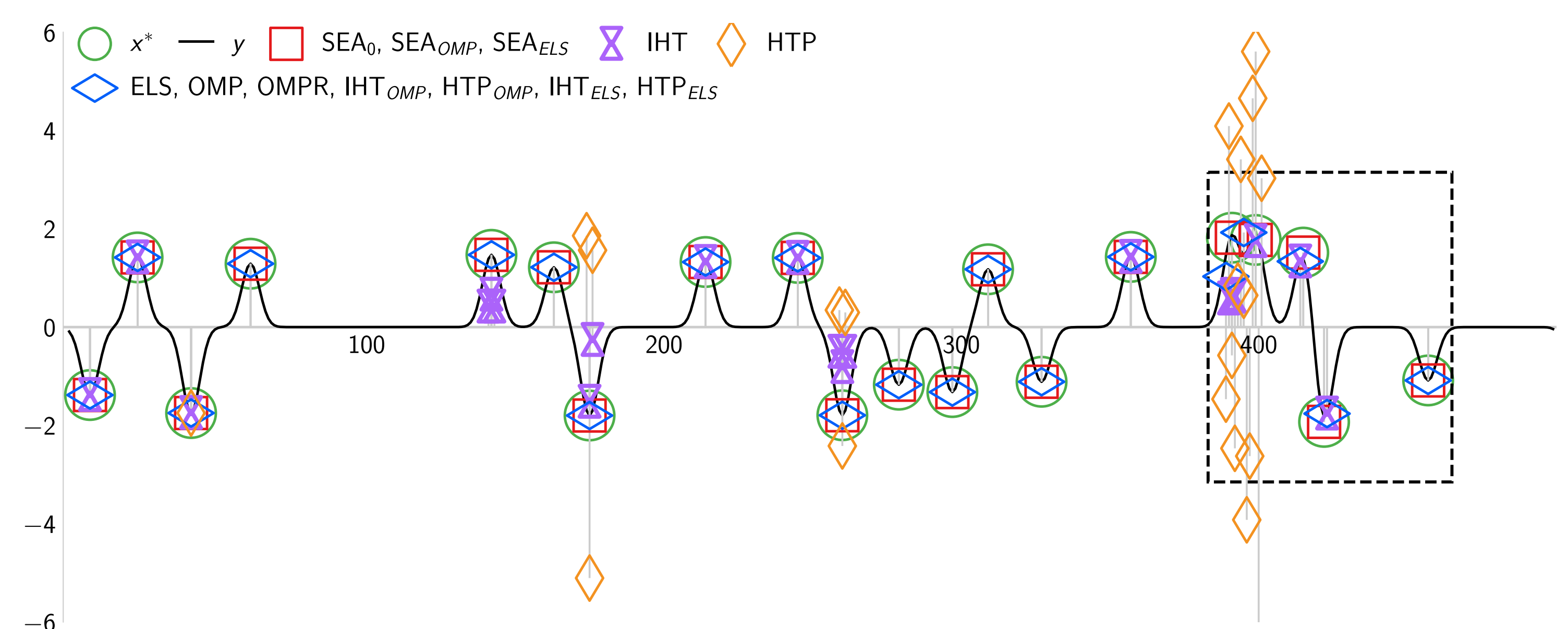


Figure 1. x^* and y with the solutions provided by the algorithms when $k = 20$

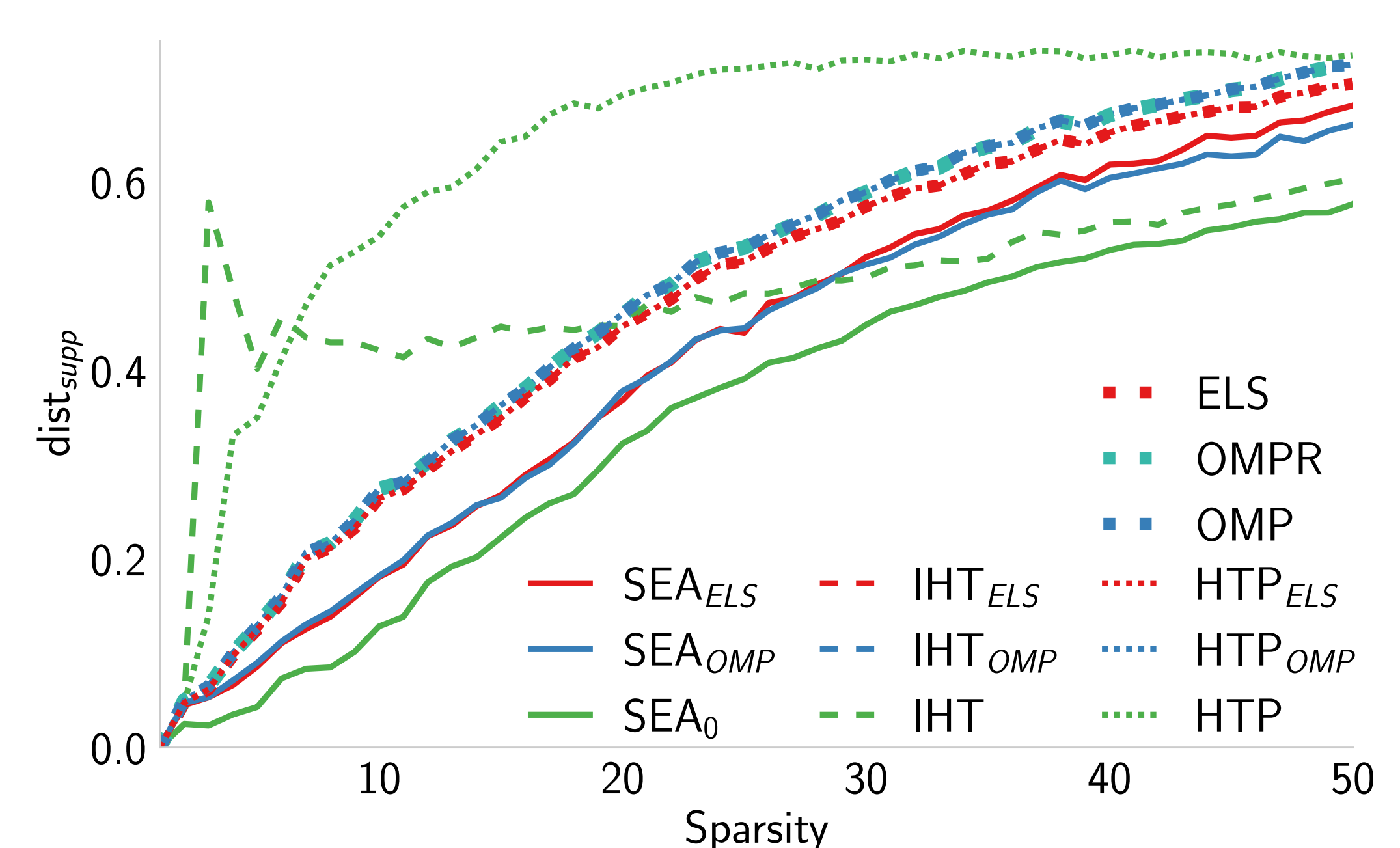


Figure 2. Average support distance $\text{supp}_{\text{dist}}(x) = \frac{k - |S^* \cap \text{supp}(x)|}{k}$ between S^* and the support of the solutions provided by the algorithms over 200 run: $\mu(A) = 0.97$, $\sigma = 3$

Phase transition diagram with $A, x_{|S^*}^* \sim \mathcal{N}(0, 1)$

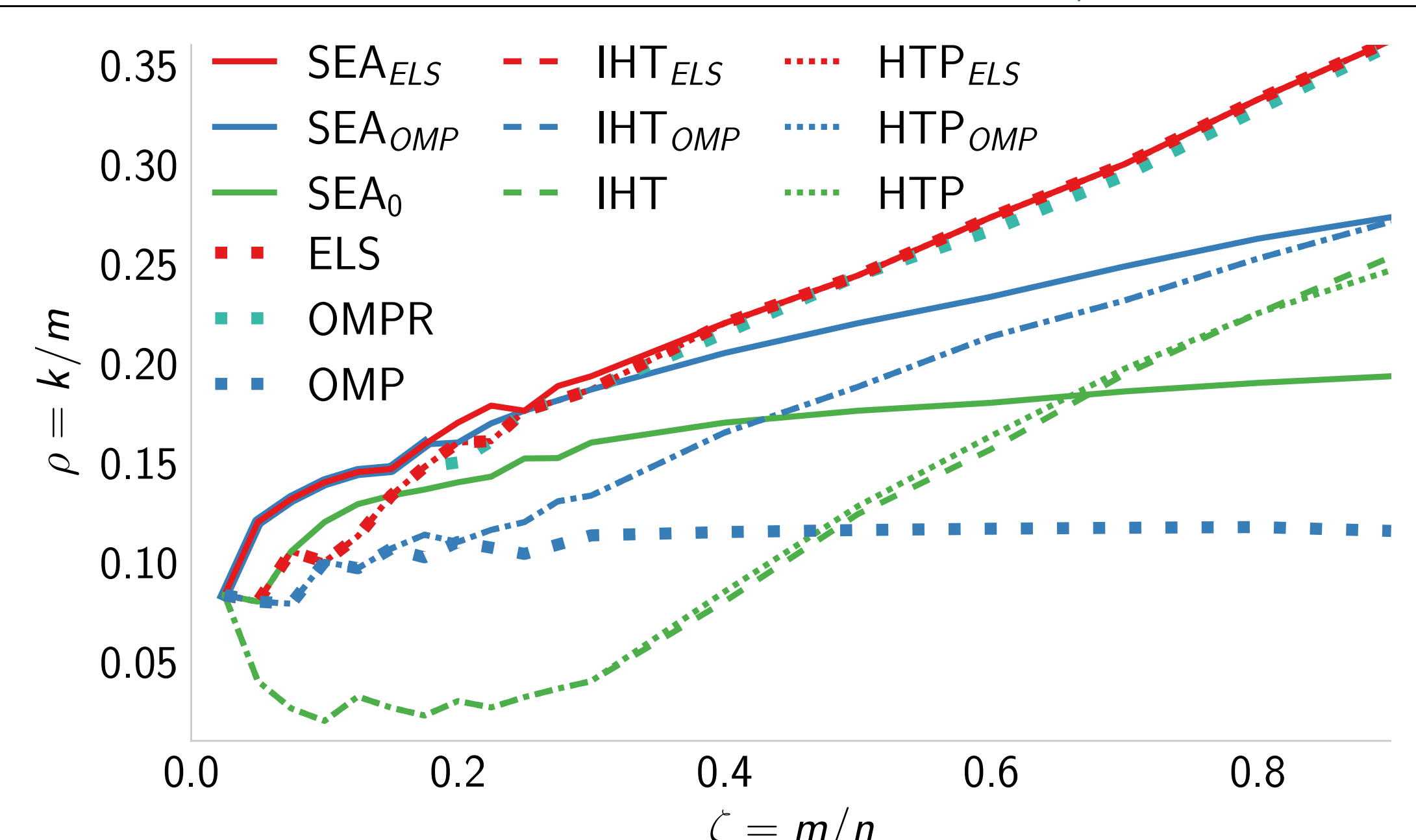


Figure 3. Empirical support recovery phase transition curves. Problems below each curve are solved by the algorithms with a success rate larger than 95% over 1000 runs