

# Large-Margin Determinantal Point Processes

Wei-Lun Chao<sup>\*1</sup>, Boqing Gong<sup>\*1</sup>, Kristen Grauman<sup>2</sup>, and Fei Sha<sup>1</sup>

<sup>1</sup>University of Southern California, <sup>2</sup>University of Texas at Austin

## Highlights

- Investigate determinantal point processes (DPPs) for discriminative subset selection
- Propose **margin** based parameter estimation to explicitly track errors in selecting subsets
- Balance different types of evaluation metrics, e.g., precision and recall
- Improve modeling flexibility by multiple-kernel based parameterization
- Attain state-of-the-art performance on the tasks of video and document summarization

## Background

- A DPP defines a probabilistic distribution over the power set of a ground set: *diverse subsets with large probabilities*

Ground set of M items,  $\mathcal{Y} = \{1, 2, \dots, M\}$

$L \in \mathbb{S}_+^M$ : a kernel matrix of pairwise similarities

$$P(y \subseteq \mathcal{Y}; L) = \frac{\det(L_y)}{\det(L + I)}$$

$$P(y = \{i, j\}; L) \propto \det(L_{\{i, j\}}) = L_{ii}L_{jj} - L_{ij}^2$$

- DPPs offer a powerful approach to modeling **diversity** in applications where the goal is to select a diverse subset from a ground set of items (e.g., retrieval, summarization)
- MAP inference (NP-hard):  $y^{\text{MAP}} = \arg\max_y P(y; L)$
- Estimate the kernel  $L$  from labeled data  $\{(y^{*(n)}, \mathcal{Y}^{(n)})\}$ 
  - Reparameterization:  $L^{(n)}(\mathcal{Y}^{(n)}; \cdot)$
  - Standard method: Maximum likelihood estimation (MLE)

$$L^{\text{MLE}} = \arg\max_L \sum_n \log P(y^{*(n)}; L^{(n)}(\mathcal{Y}^{(n)}; \cdot))$$

## Problems with existing methods

### Statistical challenges

- Limited number of training samples

### Modelling challenges

- Limited power in parameterizing kernels with the widely used quality-diversity (QD) decomposition
- Unable to track discriminative errors in selecting subsets
- Unable to differentiate different types of metrics (e.g., precision, recall) for complex structured prediction tasks

## Contribution I: Multiple-kernel representation (MKR)

- Quality-diversity (QD) decomposition

$$\forall i \in \mathcal{Y} \begin{cases} x_i: \text{quality features} \\ w_i: \text{similarity features} \end{cases} \Rightarrow \begin{cases} L_{ij} = q_i q_j S_{ij} = q_i q_j w_i^T w_j \\ q_i = q(x_i) = \exp(\cdot^T x_i) \end{cases}$$

- Multiple-kernel representation (MKR)**

$$S_{ij} = \sum_k r_k \exp\left\{-\frac{\|w_i - w_j\|_2^2}{t_k^2}\right\} + S w_i^T w_j, \quad \text{s.t.} \quad \sum_k r_k + S = 1$$

### Selected references:

- [1] A. Kulesza and B. Taskar. Determinantal point processes for machine learning. 2012.
- [2] H. T. Dang. Overview of DUC 2005. In Document Understanding Conf., 2005.
- [3] S. E. F. de Avila et al. VSUMM: A mechanism designed to produce static video summaries and a novel evaluation method. Pattern Recognition Letters, 2011.

**Acknowledgements:** F. S., B. G., and W. C. are supported by ARO Award #W911NF-12-1-0241, DARPA Award #D11AP00278, NSF IIS Award #1065243, ONR #N00014-12-1-0066, and Alfred P. Sloan Fellowship. K. G. is supported by ONR YIP #N00014-12-1-0754.

## Contribution II: Margin based parameter estimation (LME)

- Maintain desired margins between correct/incorrect subsets

$$\log P(y^*; L) \geq \max_{y \subseteq \mathcal{Y}} \log \{ \ell(y^*, y) P(y; L) \} \\ = \max_{y \subseteq \mathcal{Y}} \log \ell(y^*, y) + \log P(y; L)$$

➤ *The multiplicative margin leads to tractable optimization*

- Measure the subset discrepancy by structured loss functions

$$\ell_S(y^*, y) = \underbrace{\sum_{i \notin y^*} \mathbb{I}[i \in y]}_{\text{precision}} + \underbrace{\tilde{S} \sum_{i \in y^*} \mathbb{I}[i \notin y]}_{\text{recall}}$$

- Optimization: Jensen's inequality (softmax) for tractability

$$\log P(y^*; L) \geq \text{softmax}_{y \subseteq \mathcal{Y}} \log \ell_S(y^*, y) + \log P(y; L) \\ = \log \left( \sum_{i \notin y^*} K_{ii} + \tilde{S} \sum_{i \in y^*} (1 - K_{ii}) \right), \text{ where } K = L(L + I)^{-1}$$

- Objective function: hinge loss  $[\cdot]_+ = \max(0, \cdot)$

$$\min \sum_n \left[ -\log P(y^{*(n)}; L^{(n)}) + \log \left( \sum_{i \notin y^{*(n)}} K_{ii}^{(n)} + \tilde{S} \sum_{i \in y^{*(n)}} (1 - K_{ii}^{(n)}) \right) \right]_+$$

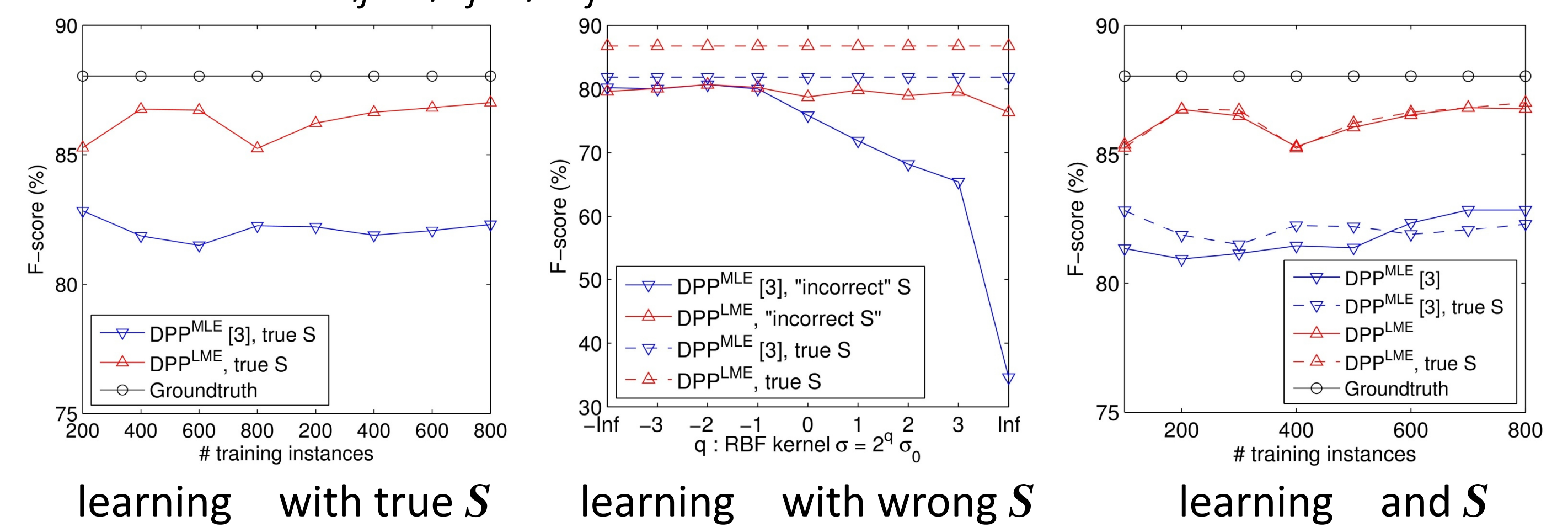
## Experiments

**Evaluation:** Precision, Recall, F-score (harmonic mean of P, R)

**Inference:** Brute-force search, minimum Bayesian risk (MBR)

### Synthetic dataset:

- Generate ground sets and target subsets by sampling  $\{x_i = \cdot\}$ , computing  $S_{ij} = q_i q_j$ , brute-force search, and adding noise



- Our method (**MKR+LME**) is more robust in parameter estimation

### Document summarization: Document Understanding Conf. (DUC)

- Train: DUC 2003 (60 clusters) Test: DUC 2004 (50 clusters)

Method	ROUGE1-F	ROUGE1-P	ROUGE1-R	ROUGE2-F	ROUGE2-P	ROUGE2-R
PEER65 [2]	37.9	37.6	38.2	9.13	--	--
DPP <sup>MLE</sup> +COS	37.9±0.08	37.4±0.08	38.5±0.08	7.72±0.06	7.63±0.06	7.83±0.06
Ours (DPP <sup>LME</sup> +COS)	38.4±0.09	37.7±0.10	39.1±0.08	8.20±0.07	8.07±0.07	8.35±0.07
Ours (DPP <sup>MLE</sup> +MKR)	39.1±0.08	39.0±0.09	39.3±0.09	9.25±0.08	9.24±0.08	9.27±0.08
Ours (DPP <sup>LME</sup> +MKR)	<b>39.7±0.05</b>	<b>39.6±0.08</b>	<b>39.9±0.06</b>	<b>9.40±0.08</b>	<b>9.38±0.08</b>	<b>9.43±0.08</b>

### Video summarization: Open Video Project (OVP)

- 50 videos from OVP, 5 user annotations, 5-fold cross-validation

Metric	VSUMM[3]		Dpp <sup>MLE</sup>	Ours (DPP <sup>LME</sup> +MKR)		
	Type1	Type2	+MKR	$\omega=2^{-6}$	$\omega=2^0$	$\omega=2^6$
F	70.25	68.20	72.94	71.25	<b>73.46</b>	72.39
P	70.57	73.14	68.40	<b>74.00</b>	69.68	67.19
R	75.77	69.14	82.51	72.71	81.39	<b>83.24</b>



**Our method can balance different types of evaluation metrics and achieve better performance in both summarization tasks**