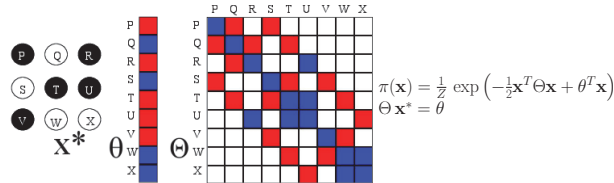


Fast, Exact & Multi-Scale Inference for Semantic Image Segmentation with Deep Gaussian CRFs

Siddhartha Chandra & Iasonas Kokkinos
INRIA GALEN & Centrale Supélec Paris



Gaussian Conditional Random Fields



Unique and exact global optimum, pairwise interactions discovered from the data via end-to-end deep learning, and fast inference via efficient implementation.

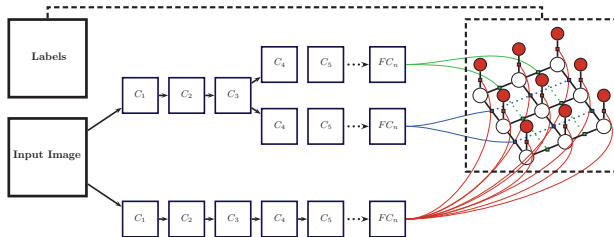
Quadratic Energy Optimization

$$E(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T (A + \lambda \mathbf{I}) \mathbf{x} - B \mathbf{x} \quad (1)$$

If $(A + \lambda \mathbf{I})$ is symmetric positive definite, **unique global minimum** at,
 $(A + \lambda \mathbf{I}) \mathbf{x} = B.$ (2)

Inference involves solving a **system of linear equations**.

Quadratic Optimization in Deep Learning



- Network populates unary and pairwise terms
- QO module proposes scores after inference
- Model parameters learnt end-to-end for arbitrary global loss (objective) \mathcal{L}
- Gradient expressions:**
 - $(A + \lambda \mathbf{I}) \frac{\partial \mathcal{L}}{\partial B} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}}, \quad (3) \quad \frac{\partial \mathcal{L}}{\partial A} = -\frac{\partial \mathcal{L}}{\partial B} \otimes \mathbf{x}. \quad (4)$
- Gradient computed **analytically** by solving a **system of linear equations**

Potts Type Model with Shared Pairwise Terms

Notation: $\hat{A}_{p_i, p_j}(l_i, l_j)$ is the pairwise energy term for pixels p_i, p_j taking the labels l_i, l_j . Per-class scores and unaries are denoted by \mathbf{x}_k , and \mathbf{b}_k , where $k \in \{1, \dots, L\}$.

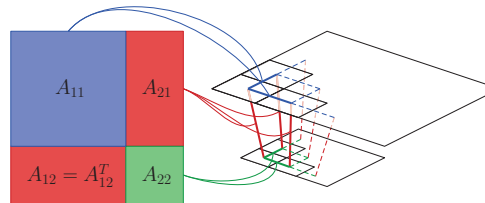
$$A_{p_i, p_j}(l_i, l_j) = \begin{cases} 0 & l_i = l_j \\ A_{p_i, p_j} & l_i \neq l_j. \end{cases} \quad (5)$$

- Fewer parameters ($P \times P$ terms) compared to general setting ($PL \times PL$ terms) for P pixels, L labels. Reduction factor of **441** for VOC Pascal
- Algebraic simplifications enable us to infer scores for each label independently

$$(\lambda \mathbf{I} + (L - 1) \hat{A}) \sum_i \mathbf{x}_i = \sum_i \mathbf{b}_i, \quad (6) \quad (\lambda \mathbf{I} - \hat{A}) \mathbf{x}_k = \mathbf{b}_k - \hat{A} \sum_i \mathbf{x}_i. \quad (7)$$

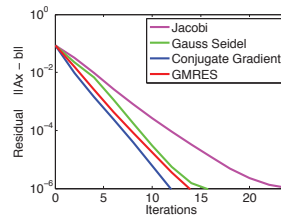
- Training **3x** faster, inference **6x** faster than general setting on VOC Pascal

Multi-Resolution Architecture



- Model parameters can capture pairwise terms between pixels across scales
- Information flow across scales
- Two kinds of interactions
 - Pairwise constraints between pixels at each resolution (Blue and Green)
 - Pairwise constraints between the same image region at different resolutions
- Inter-resolution constraints encourage pixels to share labels across resolutions

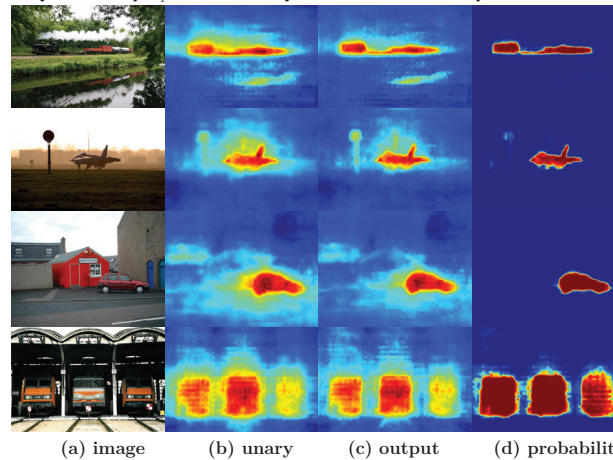
Implementation Details and Efficiency



- Conjugate Gradient > other algorithms
- Caffe based implementation using efficient CUDA Sparse, Blas routines
- General Inference time $\sim 0.02s$
- Potts-type inference time $\sim 0.003s$
- Code available at <https://github.com/siddharthachandra/gcrf>.

Experimental Setup

- All methods use *VOC PASCAL 2012* image segmentation benchmark
- Basenet is a 3-resolution variant of **Deeplab-LargeFOV**
- We experiment with 4 variants of our method
 - QO**: General pairwise terms **QO***: Potts-type *shared* pairwise terms
 - QO^{res}**: One QO per resolution **QO^{res}***: Multi-resolution QO



Results

- Comparison with Basenet:

| Method | IoU (%) | IoU after Dense CRF (%) |
|-------------------|---------|-------------------------|
| Basenet | 72.72 | 73.78 |
| QO | 73.41 | 75.13 |
| QO* | 73.20 | 75.41 |
| QO ^{res} | 73.86 | 75.46 |

- Comparison with *previously published* approaches:

| Method | mean IoU (%) |
|---|-----------------|
| Deeplab-Cross-Joint (Chen et al. ICCV 2015) | 73.9 |
| CRFRNN (Zheng et al. ICCV 2015) | 74.7 |
| Basenet | 73.8 |
| QO | 75.1 |
| QO* | 75.4 |
| QO ^{res} | 75.5 |
| Deeplab-V2 (Chen et al. Arxiv 2016) | <i>new</i> 79.7 |
| QO* + Deeplab-V2 | <i>new</i> 80.2 |



Acknowledgements: This work has been funded by the EU Projects MOBOT FP7-ICT-2011-600796 and I-SUPPORT 643666 #2020.