

# FGA-NAS: Fast Resource-Constrained Architecture Search by Greedy-ADMM Algorithm

---

Yifei Chen<sup>\*†</sup>, Junge Zhang<sup>\*†</sup>, Qiaozhe Li<sup>†</sup>, Hao Chen<sup>\*†</sup>, Kaiqi Huang<sup>\*†</sup>

*\*School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing, China*

*†CRISE, Institute of Automation, Chinese Academy of Sciences, Beijing, China*

chenyifei2019@ia.ac.cn, jgzhang@nlpr.ia.ac.cn, liqiaozhe2015@ia.ac.cn,

chenhao2019@ia.ac.cn, kqhuang@nlpr.ia.ac.cn



# CONTENTS

CONTENTS

PART 01 Motivation

PART 02 Methodology

PART 03 Experiments

PART 04 Conclusion



# 1. Motivation

Bi-Level Optimization :

$$\alpha^* = \underset{\alpha}{\operatorname{argmin}} L_{\text{val}}(\mathbf{w}^*(\alpha), \alpha)$$

$$\text{s. t. } \mathbf{w}^*(\alpha) = \underset{\mathbf{w}}{\operatorname{argmin}} L_{\text{train}}(\mathbf{w}, \alpha)$$

Hessian matrix of  $\mathbf{w}$  :

$$\frac{\partial L_{\text{val}}(\mathbf{w}^*(\alpha), \alpha)}{\partial \alpha} = \frac{\partial L_{\text{val}}(\mathbf{w}^*, \alpha)}{\partial \alpha} + \left( \frac{d\mathbf{w}^*(\alpha)}{d\alpha} \right)^T \frac{\partial L_{\text{val}}(\mathbf{w}^*, \alpha)}{\partial \mathbf{w}}$$

$$\frac{d\mathbf{w}^*(\alpha)}{d\alpha} = \left( \frac{\partial^2 L_{\text{val}}(\mathbf{w}^*, \alpha)}{\partial \mathbf{w}^2} \right)^{-1} \frac{\partial^2 L_{\text{val}}(\mathbf{w}^*, \alpha)}{\partial \mathbf{w} \partial \alpha}$$



# 1. Motivation

Bi-Level Optimization :

$$\alpha^* = \underset{\alpha}{\operatorname{argmin}} L_{\text{val}}(\mathbf{w}^*(\alpha), \alpha)$$

$$\text{s. t. } \mathbf{w}^*(\alpha) = \underset{\mathbf{w}}{\operatorname{argmin}} L_{\text{train}}(\mathbf{w}, \alpha)$$

Heuristic algorithm:

$$\alpha_{t+1} = \alpha_t - \eta_{\alpha} \nabla_{\alpha} L_{\text{val}}(w_t - \xi \nabla_w L_{\text{train}}(w_t, \alpha_t), \alpha_t)$$

$$w_{t+1} = w_t - \eta_w \nabla_w L_{\text{train}}(w_t, \alpha_{t+1})$$

Performance Collapse



## 2. Methodology

Single-Level Optimization :

$$(\mathbf{w}^*, \boldsymbol{\alpha}^*) = \underset{(\mathbf{w}, \boldsymbol{\alpha})}{\operatorname{argmin}} L_{\text{train}}(\mathbf{w}, \boldsymbol{\alpha})$$

$$s.t. r_c \geq \sum_i h_{c,i} \cdot \|\alpha_i\|_0, \quad 1 \leq c \leq n_c$$

By introducing  $\boldsymbol{\beta}_c$  and  $f_c(\boldsymbol{\beta}_c)$  :

$$(\mathbf{w}^*, \boldsymbol{\alpha}^*) = \underset{(\mathbf{w}, \boldsymbol{\alpha})}{\operatorname{argmin}} L_{\text{train}}(\mathbf{w}, \boldsymbol{\alpha}) + \sum_{c=1}^{n_c} f_c(\boldsymbol{\beta}_c)$$

$$s.t. \boldsymbol{\alpha} = \boldsymbol{\beta}_c \quad 1 \leq c \leq n_c$$

$$S_c = \{\boldsymbol{\alpha} \mid r_c \geq \sum_i h_{c,i} \cdot \|\alpha_i\|_0\}, \quad f_c(\boldsymbol{\beta}_c) = \begin{cases} 0 & , \quad \boldsymbol{\beta}_c \in S_c \\ +\infty & , \quad \text{otherwise} \end{cases}$$



## 2. Methodology

Single-Level Optimization :

$$(\mathbf{w}^*, \boldsymbol{\alpha}^*) = \underset{(\mathbf{w}, \boldsymbol{\alpha})}{\operatorname{argmin}} L_{train}(\mathbf{w}, \boldsymbol{\alpha}) + \sum_{c=1}^{n_c} f_c(\boldsymbol{\beta}_c)$$

$$s. t. \quad \boldsymbol{\alpha} = \boldsymbol{\beta}_c \quad 1 \leq c \leq n_c$$

Augmented Lagrangian  $F_\rho(\omega, \beta, \alpha, m)$  :

$$F_\rho(\omega, \beta, \alpha, m) = L_{train}(\omega, \alpha) - \sum_{c=1}^{n_c} \frac{\|m_c\|_F^2}{2\rho} \\ + \sum_{c=1}^{n_c} \left[ f_c(\beta_c) + \frac{\rho}{2} \|\alpha - \beta_c\|_F^2 + \frac{m_c}{\rho} \right]$$



## 2. Methodology

Augmented Lagrangian  $F_\rho(\omega, \beta, \alpha, \mathbf{m})$  :

$$F_\rho(\omega, \beta, \alpha, \mathbf{m}) = L_{train}(\omega, \alpha) - \sum_{c=1}^{n_c} \frac{\|\mathbf{m}_c\|_F^2}{2\rho} + \sum_{c=1}^{n_c} \left[ f_c(\beta_c) + \frac{\rho}{2} \left\| \alpha - \beta_c + \frac{\mathbf{m}_c}{\rho} \right\|_F^2 \right]$$

By introducing ADMM algorithm :

$$\begin{cases} \beta_c^{t+1} = \underset{\beta}{\operatorname{argmin}} F_\rho(w^t, \alpha^t, \beta, m^t) \\ (w^{t+1}, \alpha^{t+1}) = \underset{w, \alpha}{\operatorname{argmin}} F_\rho(w^t, \alpha^t, \beta, m^t) \\ m_c^{t+1} = m_c^t + \rho(\alpha^{t+1} - \beta_c^{t+1}) \end{cases}$$



## 2. Methodology

Augmented Lagrangian  $F_\rho(\omega, \beta, \alpha, \mathbf{m})$  :

$$F_\rho(\mathbf{w}, \beta, \alpha, \mathbf{m}) = L_{train}(\mathbf{w}, \alpha) - \sum_{c=1}^{n_c} \frac{\|\mathbf{m}_c\|_F^2}{2\rho} + \sum_{c=1}^{n_c} \left[ f_c(\beta_c) + \frac{\rho}{2} \left\| \alpha - \beta_c + \frac{\mathbf{m}_c}{\rho} \right\|_F^2 \right]$$

Sub problem 1  $\beta_c^{t+1} = \underset{\beta}{\operatorname{argmin}} F_\rho(\omega^t, \alpha^t, \beta, m^t) \rightarrow 0-1$  programming :

$$\beta_c^{t+1} = \theta_c^{t+1} \odot \left( \alpha^t + \frac{\mathbf{m}_c^t}{\rho} \right)$$

$$\theta_c^{t+1} = \underset{\theta_c}{\operatorname{argmin}} \sum_i \theta_{c,i} \cdot \left( \alpha_i^t + \frac{m_{c,i}^t}{\rho} \right)^2$$

$$s. t. \quad \theta_{c,i} \in \{0,1\}, \quad r_c \geq \sum_i h_{c,i} \cdot \theta_{c,i}, \quad 1 \leq c \leq n_c$$





## 2. Methodology

Augmented Lagrangian  $F_\rho(\mathbf{w}, \boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{m})$  :

$$F_\rho(\mathbf{w}, \boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{m}) = L_{train}(\mathbf{w}, \boldsymbol{\alpha}) - \sum_{c=1}^{n_c} \frac{\|\mathbf{m}_c\|_F^2}{2\rho} + \sum_{c=1}^{n_c} \left[ f_c(\boldsymbol{\beta}_c) + \frac{\rho}{2} \left\| \boldsymbol{\alpha} - \boldsymbol{\beta}_c + \frac{\mathbf{m}_c}{\rho} \right\|_F^2 \right]$$

Sub problem 2  $(\mathbf{w}^{t+1}, \boldsymbol{\alpha}^{t+1}) = \underset{\mathbf{w}, \boldsymbol{\alpha}}{\operatorname{argmin}} F_\rho(\mathbf{w}^t, \boldsymbol{\alpha}^t, \boldsymbol{\beta}, \mathbf{m}^t)$  :

$$(\mathbf{w}^{t+1}, \boldsymbol{\alpha}^{t+1}) = \underset{\mathbf{w}, \boldsymbol{\alpha}}{\operatorname{argmin}} \left( L_{train}(\mathbf{w}, \boldsymbol{\alpha}) + \sum_{c=1}^{n_c} \frac{\rho}{2} \left\| \boldsymbol{\alpha} - \boldsymbol{\beta}_c + \frac{\mathbf{m}_c}{\rho} \right\|_F^2 \right)$$



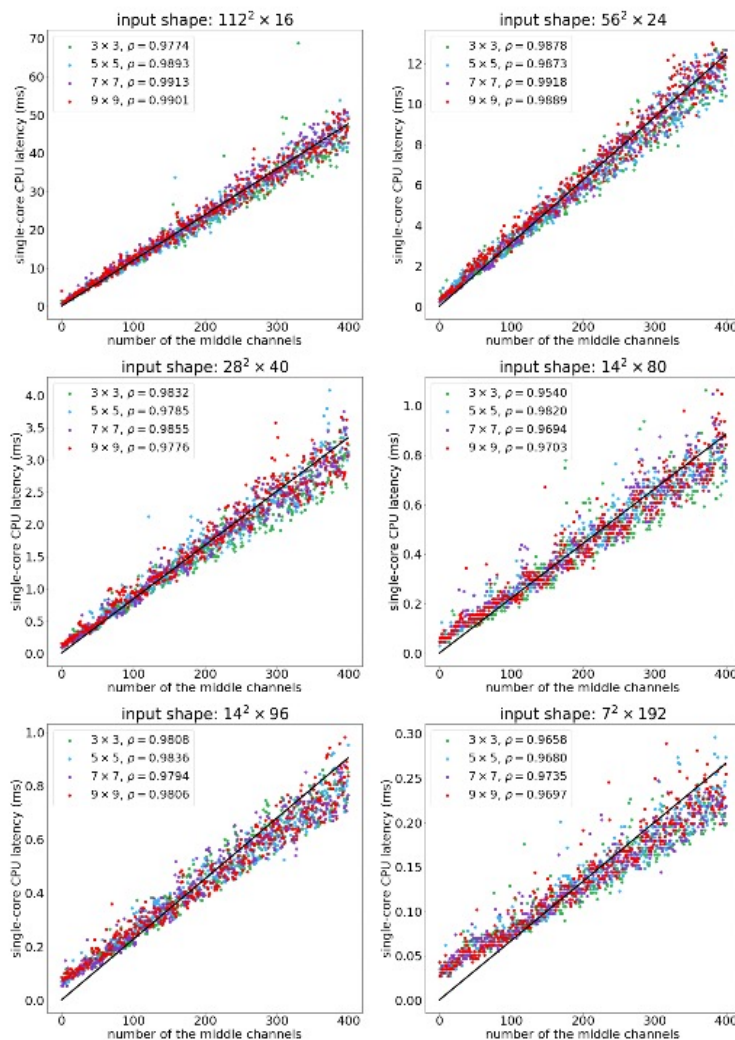
## 2. Methodology

Greedy-ADMM Search Algorithm

$$\left\{ \begin{array}{l} \beta_c^{t+1} = \underset{\beta}{\operatorname{argmin}} F_\rho(w^t, \alpha^t, \beta, m^t) \\ (w^{t+1}, \alpha^{t+1}) = \underset{w, \alpha}{\operatorname{argmin}} F_\rho(w^t, \alpha^t, \beta, m^t) \\ m_c^{t+1} = m_c^t + \rho(\alpha^{t+1} - \beta_c^{t+1}) \end{array} \right.$$

$$\left\{ \begin{array}{l} \beta_c^{t+1} = \underset{\beta}{\operatorname{argmin}} F_\rho(w^t, \alpha^t, \beta, m^t) \\ \begin{cases} w^{t+1} = w^t - \eta_w \nabla_w F_\rho(w^t, \alpha^t, \beta, m^t) \\ \alpha^{t+1} = \alpha^t - \eta_\alpha \nabla_\alpha F_\rho(w^t, \alpha^t, \beta, m^t) \end{cases} \\ m_c^{t+1} = m_c^t + \rho(\alpha^{t+1} - \beta_c^{t+1}) \end{array} \right.$$

# 3 Experiments



- (1) let the number of middle channels in the associated inverted bottleneck be one
- (2) estimate the number of parameters and the FLOPs of the above-modified inverted bottleneck
- (3) let the FLOPs and the number of parameters equal to the FLOPs and the number of parameters of the above modified block, respectively

Fig. 5. The approximation of single-core CPU latencies of inverted bottlenecks with different input shapes, different kernel sizes and different numbers of middle channels. Each point is an average of 64 measurements.

# 3 Experiments

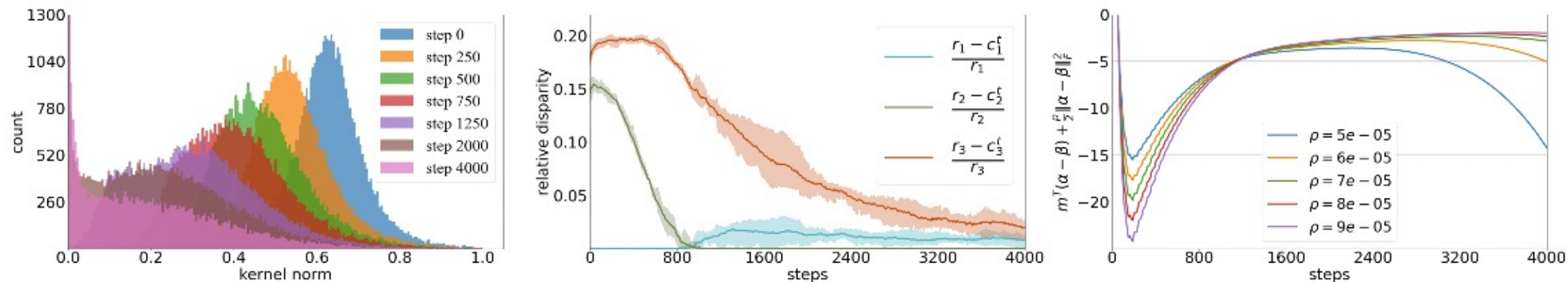


Fig. 4. **Left:** The changes in distribution of architecture parameters  $\{\alpha_{i,j}\}$  during the search. **Mid:** The change of relative disparities during the search.  $c_1^t$ ,  $c_2^t$  and  $c_3^t$  respectively denote the FLOPs, the number of parameters and the single-core CPU latency of the goal architecture at time step  $t$ .  $r_1$ ,  $r_2$  and  $r_3$  respectively denote the preset maximum FLOPs, the preset maximum the number of parameters and the preset maximum single-core CPU latency. **Right:** The convergence behavior of the search algorithm under different penalties.

# 3 Experiments



TABLE III  
COMPARISONS WITH OTHER METHODS ON IMAGENET UNDER THE MOBILE SETTING.

Model	Search			FLOPs (M)	Params (M)	Top-1 Acc (%)
	Method	Space	Cost (GPU days)			
SinglePath [15]	gradient	layer-wise	3.75	334	4.4	74.9
MobileNeXt-1.0 [32]	manual	-	-	300	3.4	74.0
MnasNet-A1 [5]	RL	stage-wise	$\geq 379$	312	3.9	75.2
HourNAS-E [13]	gradient	layer-wise	0.1	313	3.8	75.7
AtomNAS-B+ [23]	gradient	channel-wise	34	329	5.5	77.2
FBNetV2-L1 [9]	gradient	layer-wise	25	325	-	77.2
<b>FGA-NAS-A (ours)</b>	<b>gradient</b>	<b>channel-wise</b>	<b>0.2</b>	<b>320</b>	<b>5.7</b>	<b>77.0</b>
MnasNet-A2 [5]	RL	stage-wise	$\geq 379$	340	4.8	75.6
MobileNeXt-1.1 [32]	manual	-	-	420	4.28	76.7
HourNAS-F [13]	gradient	layer-wise	0.1	383	5.0	77.0
EfficientNet-B0 [31]	RL	stage-wise	$\geq 379$	390	5.3	77.3
AtomNas-C+ [23]	gradient	channel-wise	34	363	5.9	77.6
<b>FGA-NAS-B (ours)</b>	<b>gradient</b>	<b>channel-wise</b>	<b>0.2</b>	<b>358</b>	<b>6.0</b>	<b>77.4</b>



## 4 Conclusion

---

### **FGA-NAS, an efficient method for resource-constrained NAS**

- A novel search space
- A novel greedy-ADMM algorithm





中国科学院大学

University of Chinese Academy of Sciences



**Thank you.**