



首届致远学术节 学生科研成果展示

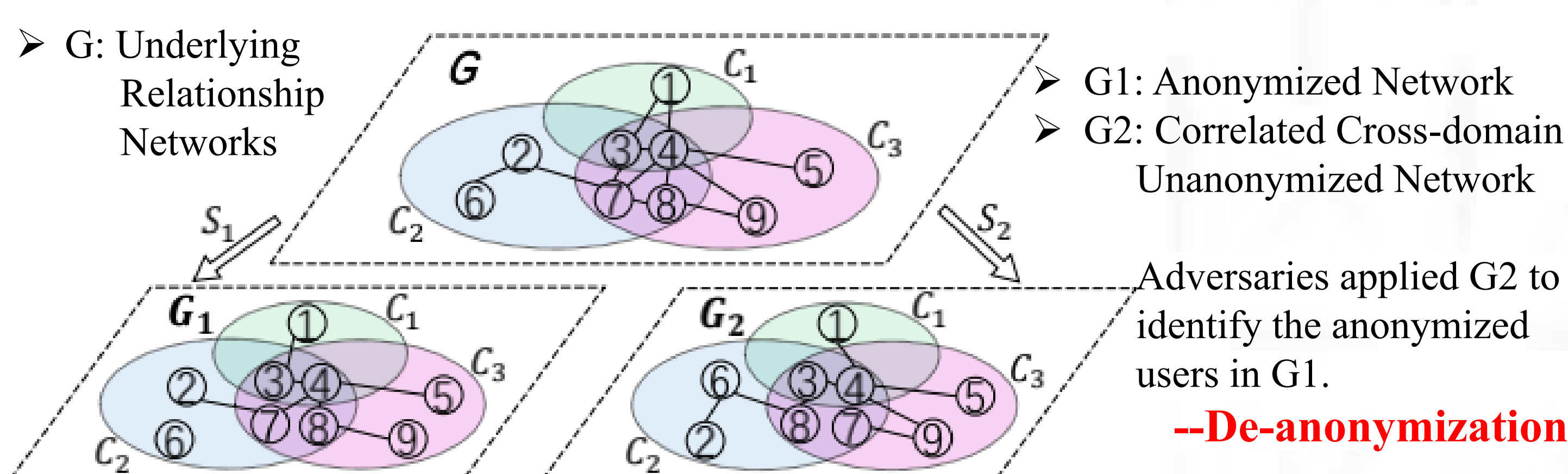
Social Network De-anonymization with Overlapping Communities: Analysis, Algorithm and Experiments

Xinyu Wu, Zhongzhao Hu, Xinzhe Fu, Luoyi Fu, Xinbing Wang, Songwu Lu

Accepted by *IEEE International Conference on Computer Communications (INFOCOM) 2018*. Acceptance Rate: 309/1606=19.2%

Background & Motivations

The advent of social networks poses severe threats on user privacy as adversaries can de-anonymize users' identities by mapping them to correlated cross-domain networks.



Goals:

- Bring a systematic analysis of this problem in **theory**, **algorithm** and **experiments**.
- Understanding the conditions under which adversaries can successfully de-anonymize users.
- Proposing efficient algorithms to solve it.
- Validating proposed algorithm under real networks.

Methodology:

1. Theory

--Modelling Overlapping Communities by the Overlapping Stochastic Block Model (OSBM).

$$Pr(C_i = \{C_{i1}, C_{i2}, \dots, C_{iQ}\}^T) = \prod_{q=1}^Q (p_q)^{C_{iq}} (1-p_q)^{1-C_{iq}} \quad Pr\{(i, j) \in E_k\} = \begin{cases} s_k & \text{if } (i, j) \in E_k \\ 0 & \text{if } (i, j) \notin E_k \end{cases}$$

--Proposing the cost function measuring de-anonymization error based on Minimum Mean Square Error (MMSE).

$$\hat{\Pi} = \arg \min_{\Pi \in \Pi^n} E_{\Pi_0} \{d(\Pi, \Pi_0)\}$$

$$= \arg \min_{\Pi \in \Pi^n} \sum_{\Pi_0 \in \Pi^n} \|\Pi - \Pi_0\|_F^2 Pr(\Pi_0 | G_1, G_2, \theta),$$

NP-hard

$$\hat{\Pi} = \arg \max_{\Pi \in \Pi^n} \sum_{\Pi_0 \in \Pi^n} \|\Pi - \Pi_0\|_F^2 \|\Pi_0 \hat{A} - \hat{B} \Pi_0\|_F^2$$

--Transforming minimizing MMSE to a polynomially solvable problem by the restriction of Sequency Inequality.

$$\hat{\Pi} = \arg \max_{\Pi \in \Pi^n} \sum_{\Pi_0 \in \Pi^n} \|\Pi - \Pi_0\|_F^2 \|W \circ (\Pi_0 A - B \Pi_0)\|_F^2$$

WEMP

Equivalent in expectation

2. Algorithm

$$\tilde{\Pi} = \arg \min_{\Pi \in \Pi^n} \|\Pi \hat{A} - \hat{B} \Pi\|_F^2$$

--Optimality of WEMP: Under mild conditions, solving WEMP ensures negligible de-anonymization error in large-scale networks.

$$\text{as } n \rightarrow \infty, \frac{\|\tilde{\Pi} - \Pi_0\|_F^2}{\|\Pi_0\|_F^2} \rightarrow 0. \quad \begin{matrix} \Pi_0 : \text{Ground-truth mapping} \\ \tilde{\Pi} : \text{Estimated mapping by WEMP} \end{matrix}$$

--Solvability of WEMP: We proposed a Convex-concave Based De-anonymization Algorithm (CBDA).

$$F_0(\Pi) = \|\hat{A} - \Pi \hat{B} \Pi^T\|_F^2 + \mu \|\Pi M - M\|_F^2, \quad \text{The Equivalent Original Objective}$$

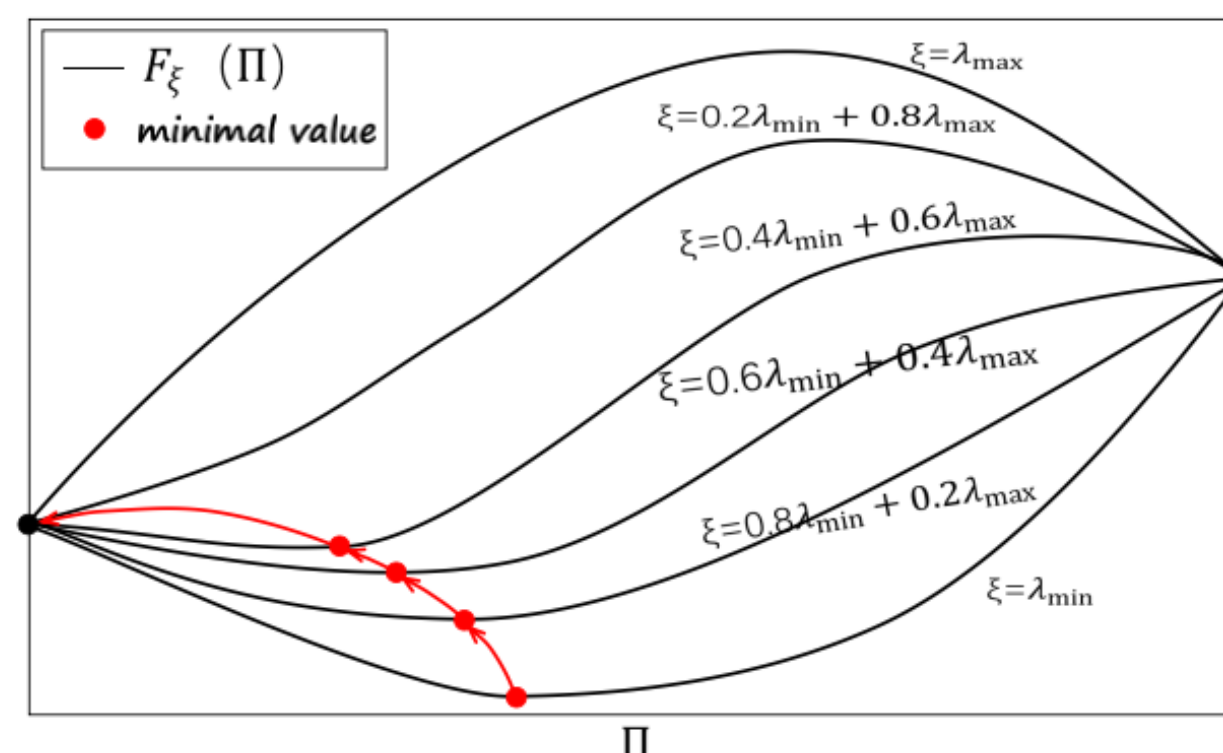
$$F_1(\Pi) = F_0(\Pi) + \frac{\lambda_{\min}}{2} (n - \|\Pi\|_F^2); \quad \text{Convex Relaxation of the Objective}$$

$$F_2(\Pi) = F_0(\Pi) + \frac{\lambda_{\max}}{2} (n - \|\Pi\|_F^2). \quad \text{Concave Relaxation of the Objective}$$

$$F(\Pi) = (1-\alpha)F_1(\Pi) + \alpha F_2(\Pi) = F_0(\Pi) + 2\xi(n - \|\Pi\|_F^2),$$

$$\xi = (1-\alpha)\lambda_{\min} + \alpha\lambda_{\max}, \quad \xi \in [\lambda_{\min}, \lambda_{\max}].$$

Combination of 2 Relaxed Versions and modification of ξ



Algorithm 1: Convex-concave Based De-anonymization Algorithm (CBDA)

Input: Adjacent matrices A and B; Community assignment matrix M;

Weight controlling parameter μ ; Adjustable parameters $\delta, \Delta\xi$.

Output: Estimated permutation matrix $\tilde{\Pi}$.

1: Form the objective function $F_0(\Pi)$ and $F(\Pi)$.

2: $\xi \leftarrow 0, k \leftarrow 1$. Initialize Π_1 . Set ξ_m , the upper limit of ξ .

3: while $\xi < \xi_m$ and $\Pi_k \notin \Omega_0$ do

4: while $k = 1$ or $|F(\Pi_{k+1}) - F(\Pi_k)| \geq \delta$ do

5: $X^\pm \leftarrow \arg \min_{X^\pm} \text{tr}(\nabla_{\Pi_k} F(\Pi_k)^T X^\pm)$, where $X^\pm \in \Omega$.

//Finding the optimal descent direction

6: $\gamma_k \leftarrow \arg \min_{\gamma} F(\Pi_k + \gamma(X^\pm - \Pi_k))$, where $\gamma_k \in [0, 1]$. //Finding the optimal step size

7: $\Pi_{k+1} \leftarrow \Pi_k + \gamma_k(X^\pm - \Pi_k), k \leftarrow k+1$.

//Estimation Update

8: end while

9: $\xi \leftarrow \xi + \Delta\xi$. Based on Frank-Wolfe algorithm

10: end while

CBDA: Overcoming the brute projection by classical convex optimization technique, which may cause large estimation error.

Experiments & Results:

Experiment Settings:

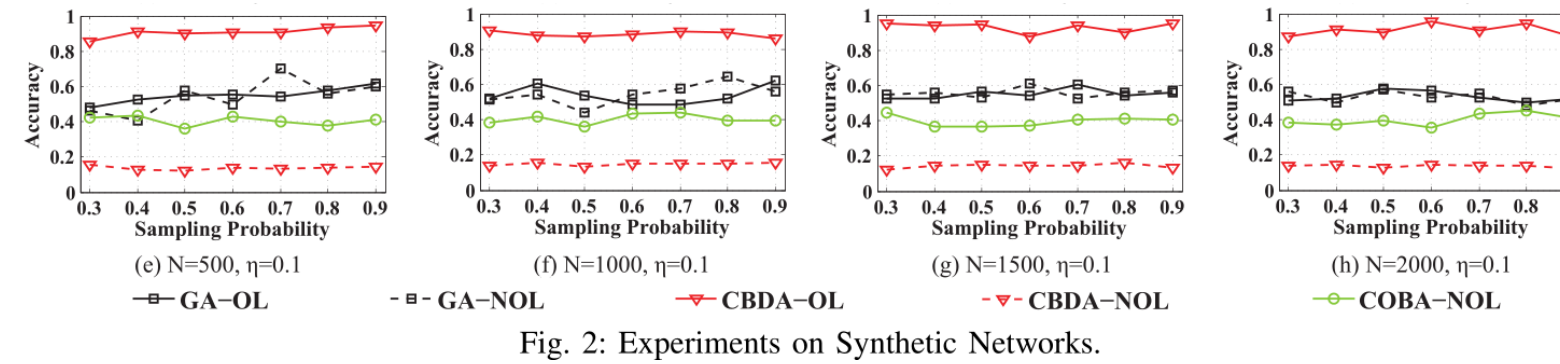
Notation	Definition	Range
N	Number of Nodes	{500, 1000, 1500, 2000}
α	Sampling Probability ($\alpha_1 = \alpha_2 = \alpha$)	{0.3, 0.4, 0.5}
η	Community Ratio	{0.05, 0.1}
OL/NOL	Overlapping or Non-Overlapping	{OL, NOL}

Experiment Datasets:

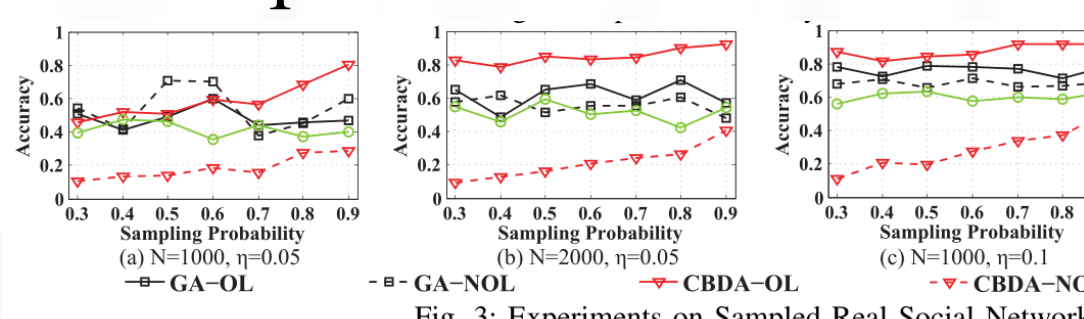
Dataset	Synthetic	Sampled Real Social	Cross-Domain Co-author
Source	OSBM	LiveJournal [1]	MAG [11]
Num. of Nodes	500 ~ 2000	500 ~ 2000	3176
Num. of Communities	25 ~ 1000	25 ~ 1000	89

Experiments Results:

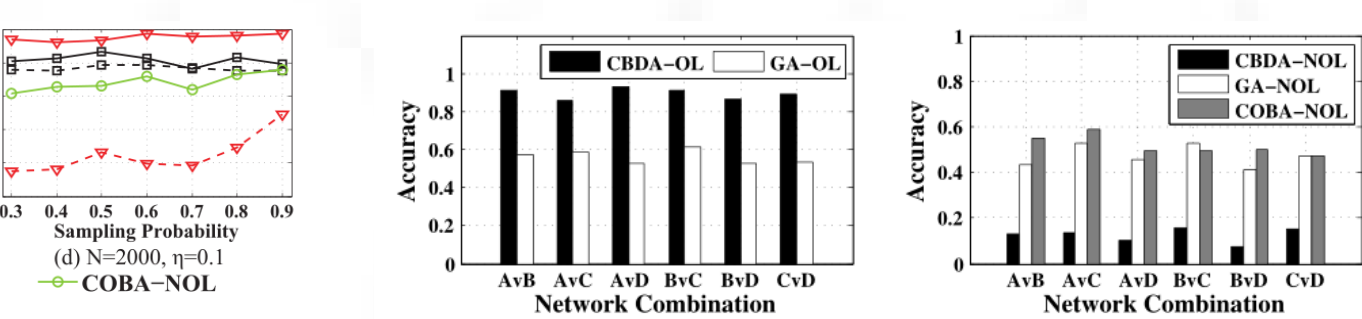
➤ Synthetic Networks.



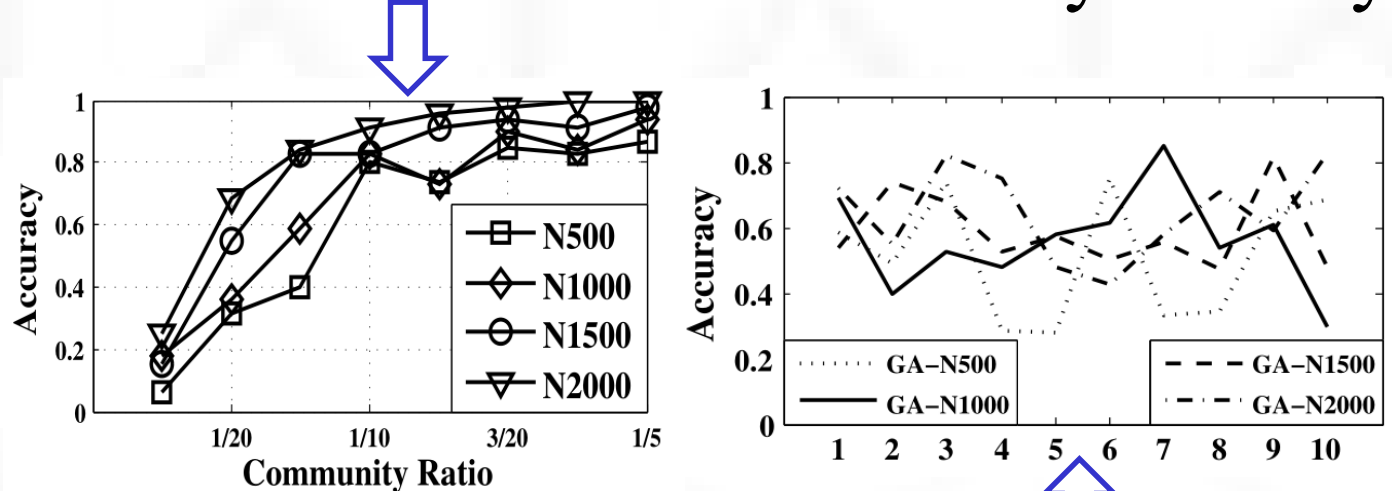
➤ Sampled Social Networks:



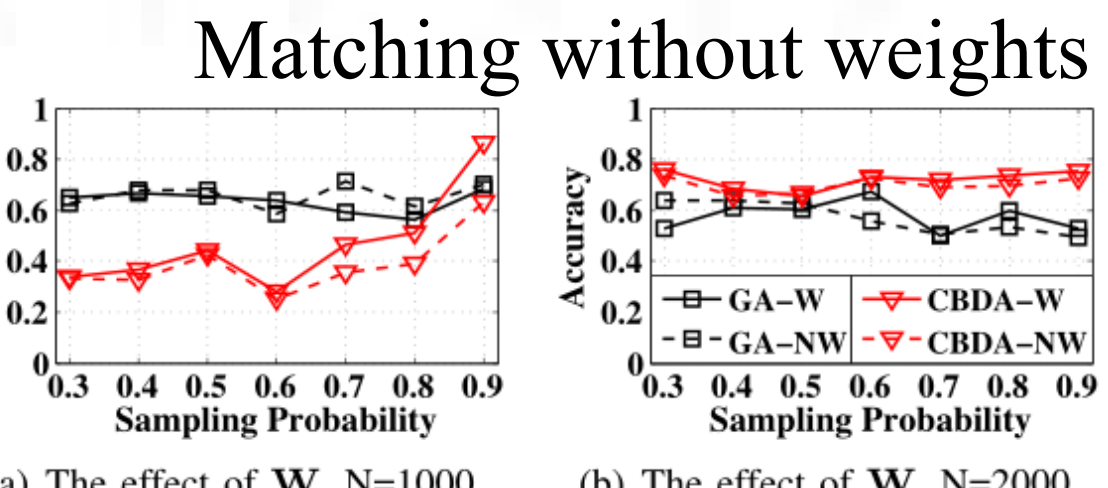
➤ Correlated Cross-Domain Networks:



➤ The Influence of Sampled Community Density



➤ MMSE vs Matching without weights



➤ The Instability of Genetic Algorithm

Conclusions:

- 1) CBDA outperforms the Genetic Algorithm and Convex Optimization Based Algorithm under networks with overlapping communities;
- 2) CBDA outperforms GA and COBA more in larger networks; (Meeting the theoretical results)
- 3) Overlapping Density positively impacts the de-anonymization accuracy; (Meeting the theoretical results)
- 4) De-anonymization based on MMSE outperforms the non-weighted cost function in prior art;
- 5) CBDA promises more practical use than GA on the factor of stability.

Paper Links:

9-page: http://wuxinyusjtu.3www.win/Infocom18_De_anonymization.pdf

Full paper: http://wuxinyusjtu.3www.win/IT_Xinyu_Wu.pdf

个人信息：吴昕宇，2014级，工科
邮箱：wuxinyu@sjtu.edu.cn