# Mathematical Statistics Team

*Shimodaira Lab*
*Graduate School of Informatics, Kyoto University*

Hidetoshi Shimodaira (Team Leader @ Kyoto Univ), Kazuki Fukui (RA), Akifumi Okuno (RA), Tetsuya Hada (RA), Masaaki Inoue (RA), Geewook Kim (RA), Thong Pham (Postdoctoral Researcher @ AIP), Tomoharu Iwata (Visiting Scientist @ NTT), Yoshikazu Terada (Visiting Scientist @ Osaka Univ), Shinpei Imori (Visiting Scientist @ Hiroshima Univ), Kei Hirose (Visiting Scientist @ Kyushu Univ )

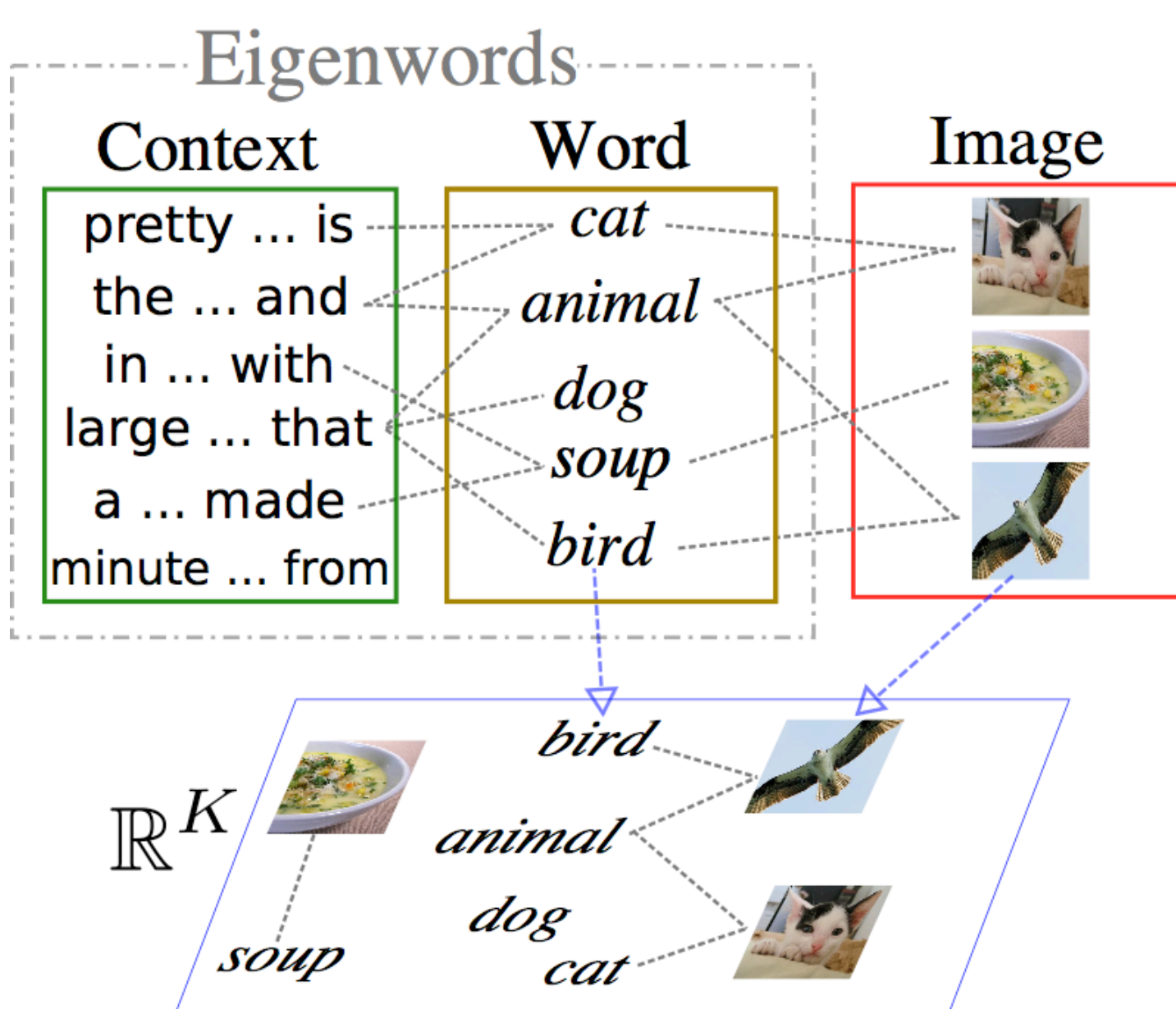## Statistics and Machine Learning: Methodology and Applications

- Inductive inference, Resampling methods, Information geometry
- Generalization error under Missing, Covariate-shift, etc.
- Multi-modal data representation, Graph Embedding, and Multivariate Analysis
- Network growth mechanism (Preferential attachment and fitness)
- Phylogenetics, Gene expression (hierarchical clustering)
- Image, word embedding (search and reasoning)

## Multimodal Eigenwords (Fukui, Oshikiri and Shimodaira, Textgraphs 2017)

- A multimodal word embedding that jointly embeds words and corresponding visual information

- We employed **Cross-Domain Matching Correlation Analysis** (CDMCA; Shimodaira 2016) for extending a CCA-based word embedding (Dhillon et al. 2015) to deal with complex associations

### Our proposed method:

- Feature learning via graph embedding
- Feature vectors reflect both semantic and visual similarities
- The method enables **multimodal vector arithmetic** between images and words

Eigenwords

Context　Word　Image

pretty … is
the … with
in … large
large … that
a … made
minute … from

*cat*
*animal*
*dog*
*soup*
*bird*

$\mathbb{R}^K$

*bird*
*animal*
*dog*
*cat*
*soup*

### Multimodal vector arithmetic:

- "day" + "night" =

- "brown" + "white" =

## Probabilistic Multi-view Graph Embedding (PMvGE)

(Okuno, Hada and Shimodaira, arXiv:1802.04630)

Multi-view feature learning with many-to-many associations via neural networks for predicting new associations

$$w_{ij}^{(de)} \sim \mathrm{Po}\left(\alpha^{(de)} \exp\left(\langle f_{\psi}^{(d)}(x_i^{(d)}), f_{\psi}^{(e)}(x_j^{(e)})\rangle\right)\right)$$

Strength of association ($\geqq 0$)　Neural network　data vector (in view-$d$)

View-1 $\mathbb{R}^{p_1}$　View-2 $\mathbb{R}^{p_2}$　View-D $\mathbb{R}^{p_D}$

Neural networks $f_{\psi}^{(d)} : \mathbb{R}^{p_d} \to \mathbb{R}^K$
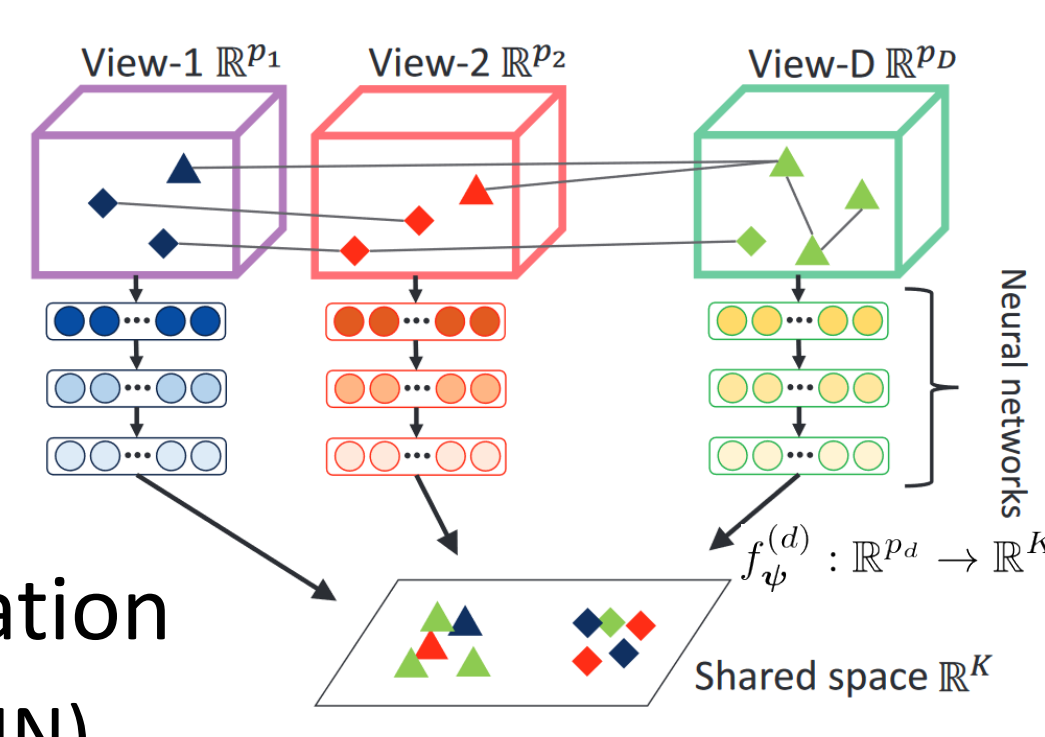
Shared space $\mathbb{R}^K$

### Theorem (universal approximation):

Inner product + Neural networks

≈ Arbitrary similarity + Arbitrary transformation

(Mercer's theorem + universal approximation of NN)

### Advantages:

(1) PMvGE with neural networks is proved to be **highly expressive**.
(2) PMvGE **approximately non-linearly generalizes CDMCA** (Shimodaira 2016, Neural Networks) which already generalizes various existing methods such as CCA, LPP, and Spectral graph embedding.
(3) Likelihood-based estimation of neural networks is **efficiently computed** by mini-batch Stochastic Gradient Descent.
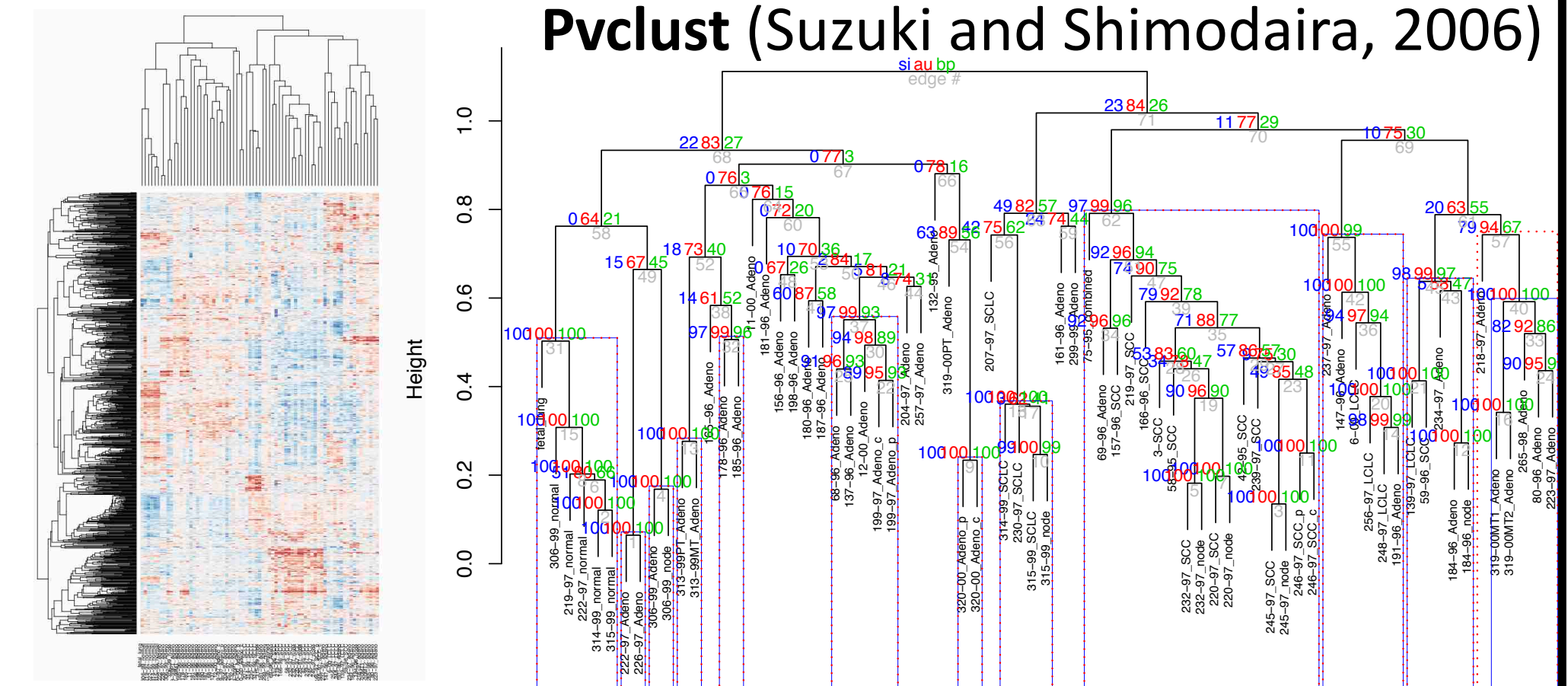
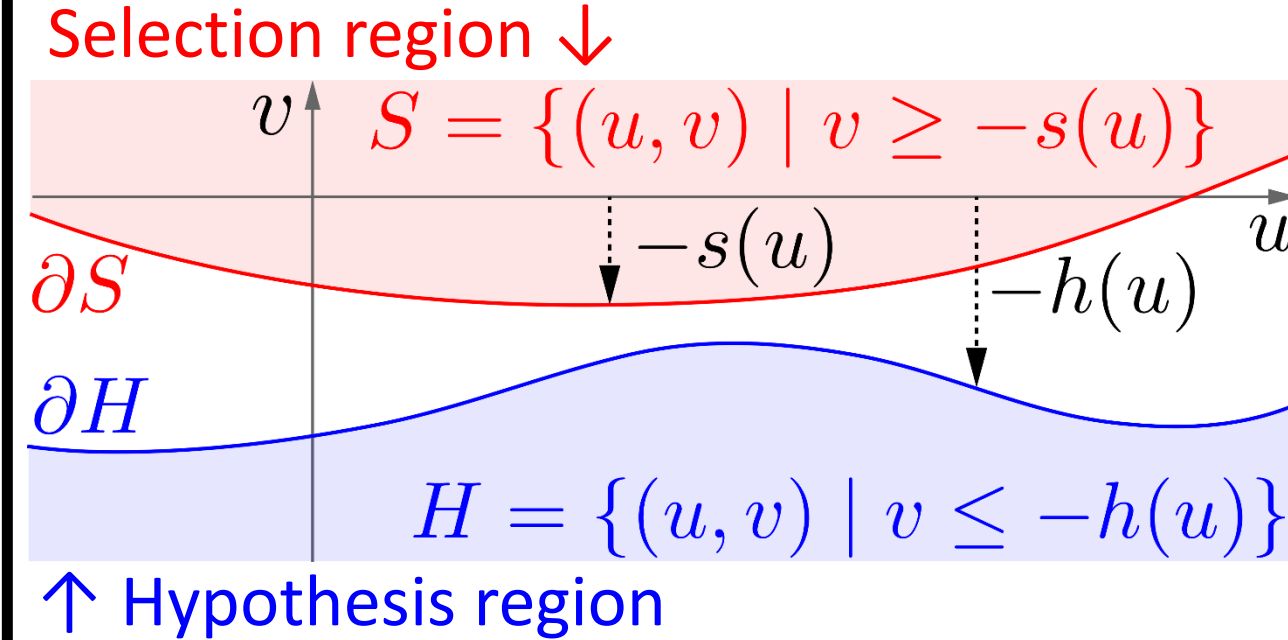## Selective inference (Terada and Shimodaira, arXiv:1711.00949)

- **Motivation** : **Assessing the confidence of each obtained cluster**
- Consider approx. unbiased p-values as frequentist confidence measures
  **Null hypothesis = obtained cluster is NOT true**
- Example : lung data (73 tumors, 916 genes; Garber et al., 2001)

**Pvclust** (Suzuki and Shimodaira, 2006)

The hypothesis is tested **only for** the clusters appeared in the obtained tree

Need **selective inference** for pvclust

- Two asymptotic theories

Selection region ↓

$v$　$S = \{(u,v) \mid v \geq -s(u)\}$　$\bar{u}$

$\partial S$　$-s(u)$　$-h(u)$

$\partial H$

$H = \{(u,v) \mid v \leq -h(u)\}$

↑ Hypothesis region

(1) Large sample theory with "**nearly parallel surfaces**"

$h(u) = h_0 + h_i u_i + h_{ij} u_i u_j + h_{ijk} u_i u_j u_k + \cdots,\ n \to \infty$

$h_0 = O(1),\ h_i = O(n^{-1/2}),\ h_{ij} = O(n^{-1/2}),\ h_{ijk} = O(n^{-1/2}),\cdots$ Smooth surface

$(u,v)$ are $O(\sqrt{n})$, but $h_0 = O(n^{1/2}),\ h_i = O(1)$ are multiplied by $O(n^{-1/2})$

(2) Asymptotic theory of "**nearly flat surfaces**"

$\sup_{u \in \mathbb{R}^m} |h(u)| = O(\lambda),\ \int_{\mathbb{R}^m} |h(u)|\, du < \infty,\ \int_{\mathbb{R}^m} |\mathcal{F}h(\omega)|\, d\omega < \infty$

(Shimodaira 2008)　non-smooth surface　$\lambda \to 0$

### Algorithm : Computing approx. unbiased selective p-values

1. Specify several $n' \in \mathbb{N}$ values, and set $\sigma^2 = n/n'$. Set the number of bootstrap replicates $B$, say, 1000.
2. For each $n'$, perform bootstrap resampling to generate $Y^*$ for $B$ times and compute $\alpha_{\sigma^2}(H|y) = C_H/B$ and $\alpha_{\sigma^2}(S|y) = C_S/B$ by counting the frequencies $C_H = \#\{Y^* \in H\}$ and $C_S = \#\{Y^* \in S\}$. (We actually work on $\mathcal{X}_{n'}^*$ instead of $Y^*$.) Compute $\psi_{\sigma^2}(H|y) = \sigma\Phi^{-1}(\alpha_{\sigma^2}(H|y))$ and $\psi_{\sigma^2}(S|y) = \sigma\Phi^{-1}(\alpha_{\sigma^2}(S|y))$.
3. Estimate parameters $\beta_H(y)$ and $\beta_S(y)$ by fitting models　Model fitting to psi
   $$\psi_{\sigma^2}(H|y) = \varphi_H(\sigma^2|\beta_H)\quad \psi_{\sigma^2}(S|y) = \varphi_S(\sigma^2|\beta_S),$$
   respectively. The parameter estimates are denoted as $\hat{\beta}_H(y)$ and $\hat{\beta}_S(y)$. If we have several candidate models, apply above to each and choose the best model based on AIC value.
4. Approximately unbiased p-values of selective inference ($p_{\mathrm{SI}}$) and non-selective inference ($p_{\mathrm{AU}}$) are computed by one of (A) and (B).
   (A) Extrapolate $\psi_{\sigma^2}(H|y)$ and $\psi_{\sigma^2}(S|y)$ to $\sigma^2 = -1$ and 0, respectively, by
   Extrapolation　$z_H = \varphi_H(-1|\hat{\beta}_H(y))$ and $z_S = \varphi_S(0|\hat{\beta}_S(y))$,
   and then compute p-values by
   **Selective p-value**　$p_{\mathrm{SI}}(H|S,y) = \dfrac{\Phi(z_H)}{\Phi(z_H + z_S)}$ and $p_{\mathrm{AU}}(H|y) = \Phi(z_H)$.　Non-selective p-value
   (B) Specify $k \in \mathbb{N},\ \sigma_0^2, \sigma_{-1}^2 > 0$ (e.g., $k = 3$ and $\sigma_{-1}^2 = \sigma_0^2 = 1$). Extrapolate $\psi_{\sigma^2}(H|y)$ and $\psi_{\sigma^2}(S|y)$ to $\sigma^2 = -1$ and 0, respectively, by
   $z_{H,k} = \varphi_{H,k}(-1|\hat{\beta}_H(y), \sigma_{-1}^2)$ and $z_{S,k} = \varphi_{S,k}(0|\hat{\beta}_S(y), \sigma_0^2)$,
   where the Taylor polynomial approximation of $\varphi_H$ at $\tau^2 > 0$ with $k$ terms is:
   $$\varphi_{H,k}(\sigma^2|\hat{\beta}_H(y), \tau^2) = \sum_{j=0}^{k-1} \frac{(\sigma^2 - \tau^2)^j}{j!} \left.\frac{\partial^j \varphi_H(\sigma^2|\hat{\beta}_H(y))}{\partial(\sigma^2)^j}\right|_{\sigma^2 = \tau^2}$$
   and that of $\varphi_S$ is defined similarly. Then compute p-values by
   $p_{\mathrm{SI},k}(H|S,y) = \dfrac{\Phi(z_{H,k})}{\Phi(z_{H,k} + z_{S,k})}$ and $p_{\mathrm{AU},k}(H|y) = \Phi(z_{H,k})$.

### Theorem (large sample theory):

Boundary surfaces of $H$ and $S$ are **smooth** and "**nearly parallel**", ⇒ The proposed p-value is **second order accurate with error** $O(n^{-1})$

### Theorem (nearly flat surfaces):

Boundary surfaces are "**nearly flat**", approaching flat but **allowing non-smooth** such as cones and polyhedra ⇒ the proposed p-value is justified **as unbiased with error** $O(\lambda^2)$

### Key point: **Multiscale Bootstrap**
(Shimodaira, 2002; 2004)

✓ Low computational cost : $O(B)$
✓ **Double bootstrap** method has **same accuracy** but high comp. cost $O(B^2)$

For non-smooth surfaces (such as cones and polyhedra)
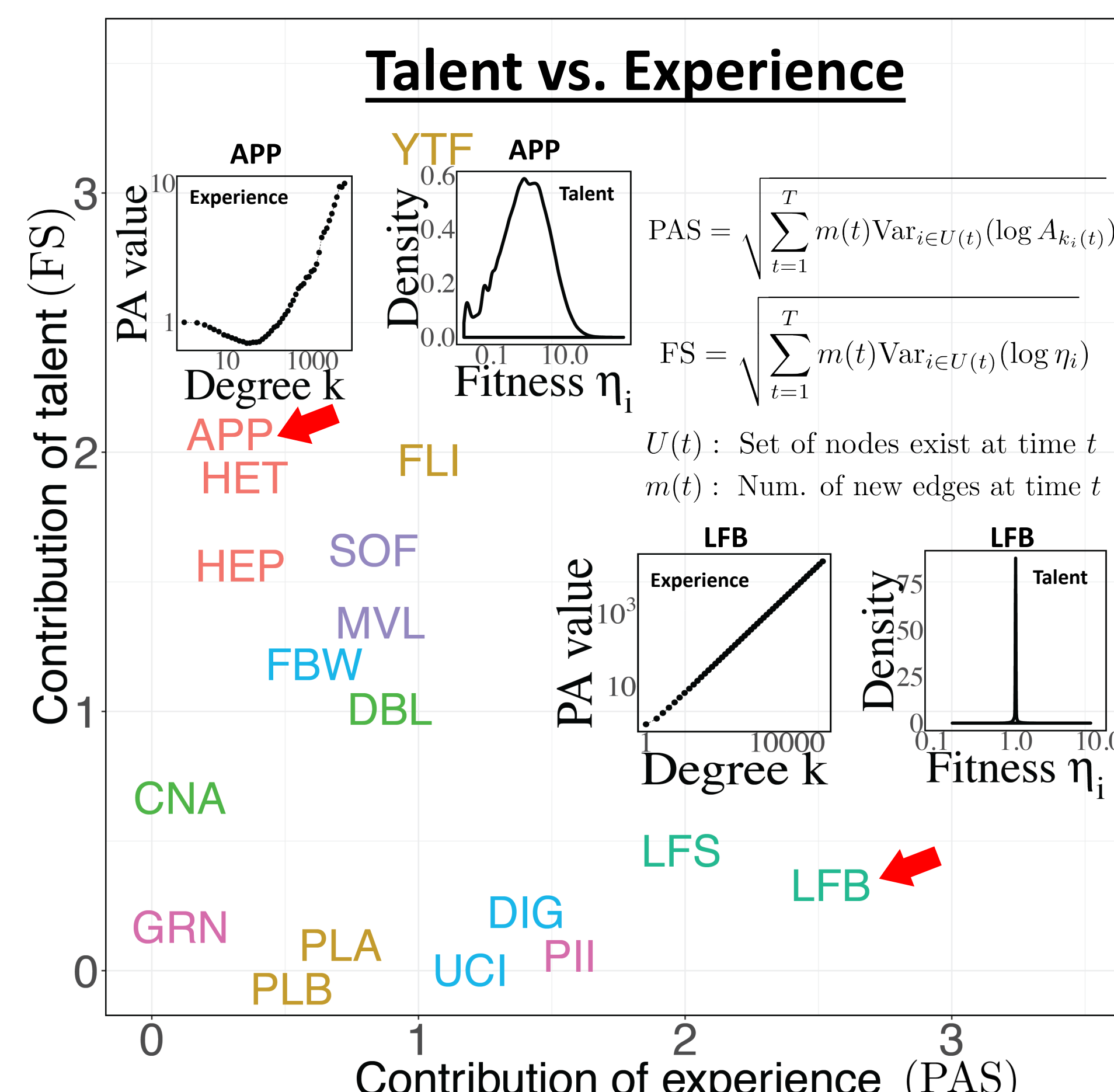
## PAFit: an R Package for Estimating Preferential Attachment and Node Fitness in Temporal Complex Networks

(Pham, Sheridan and Shimodaira, arXiv:1704.06017)

- What drives the growth of real-world networks across diverse fields?
- Using two **interpretable and universal** mechanisms: (1) **Talent**: the intrinsic ability of a node to attract connections (fitness $\eta_i$) and (2) **Experience**: the preferential attachment (PA) function $A_k$ governing the extent to which having more connections makes a node more/less attractive for forming new connections in the future.

Probability that node $v_i$ gets a new edge at time $t = A_{k_i(t)} \eta_i$

### Key finding: **Although both talent and experience contribute to the growth process, the ratios of their contributions vary greatly.**

**Talent vs. Experience**

APP　YTF　APP

PA value　Density
Experience　Talent

Degree k　Fitness $\eta_i$

$$\mathrm{PAS} = \sqrt{\sum_{t=1}^{T} m(t) \mathrm{Var}_{i \in U(t)}(\log A_{k_i(t)})}$$

$$\mathrm{FS} = \sqrt{\sum_{t=1}^{T} m(t) \mathrm{Var}_{i \in U(t)}(\log \eta_i)}$$

$U(t)$ : Set of nodes exist at time $t$
$m(t)$ : Num. of new edges at time $t$

LFB　LFB
Experience　Talent

PA value　Density
Degree k　Fitness $\eta_i$

Contribution of talent (FS)

APP
HET　FLI
HEP　SOF
MVL
FBW
DBL
CNA
GRN　DIG
PLA　PII
PLB　UCI
LFS　LFB

Contribution of experience (PAS)

1. Citation networks
HEP: arXiv hep–ph papers
HET: arXiv hep–th papers
APP: APS journal papers
2. Social networks
YTF: YouTube followship
PLB: Prosper loan before Lehman
PLA: Prosper loan after Lehman
FLI: Flickr friendship
3. Co–author networks
DBL: dblp authors
CNA: complex network authors
4. Interaction networks
LFS: last.fm song
LFB: last.fm band
5. Communication networks
UCI: UCIrvine forum message
FBW: Facebook wall–post
DIG: Digg message
6. Rating networks
SOF: Stackoverflow favourite
MVL: MovieLens rating
7. Biological networks
GRN: Cancer gene
PII: Human PPI