

Phương pháp và công nghệ dữ liệu Lớn

Machine learning for Big Data

Phùng Quốc Định

Centre for Pattern Recognition and Data Analytics

Deakin University, Australia

Email: dinh.phung@deakin.edu.au

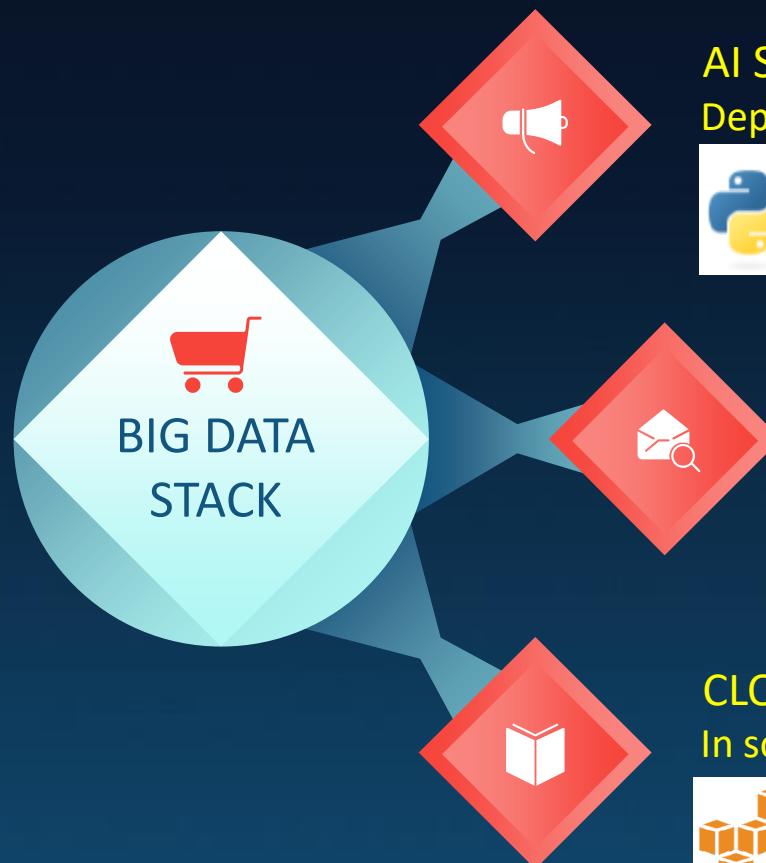
(published under Dinh Phung)

Outline

- Mô hình lớn cho dữ liệu lớn
 - Why traditional methods might fail on big data?
 - Three approaches to scale up big models
 - Scale up Gradient Descent Methods
 - Scale up Probabilistic Inference
 - Scale up Model Evaluation (bootstrap)
- Ba cách tiếp cận để phát triển mô hình lớn
 - Data augmentation
 - Stochastic Variational Inference for Graphical Models
 - Stochastic Gradient Descent and Online Learning
- Công nghệ nào cho mô hình lớn?
 - Hadoop, Spark, TensorFlow

Tiếp cận Dữ liệu lớn

- Technical approach a big data problem



AI SYSTEMS AND APPLICATIONS

Deployment into applications and products



TECHNOLOGY for DATA PROCESSING

Pick a suitable processing technology or platform



CLOUD COMPUTING INFRASTRUCTURE

In some cases, one can decide to build their own



Mô hình lớn cho dữ liệu lớn

Scale up Machine Learning
and Statistical Models

- Mô hình lớn cho dữ liệu lớn
 - Three approaches to scale up big models
 - Scale up Gradient Descent Methods
 - Scale up Probabilistic Inference
 - Scale up Model Evaluation (bootstrap)
- Ba cách tiếp cận để phát triển mô hình lớn
 - Data augmentation
 - Stochastic Variational Inference for Graphical Models
 - Stochastic Gradient Descent and Online Learning
- Công nghệ nào cho mô hình lớn?
 - Hadoop, Spark, TensorFlow

Scaling up ML and Stats models

Mô hình tính toán cho dữ liệu lớn

- Tại sao nhiều phương pháp truyền thống trong ML và Stats có thể không dùng được trong bài toán DLL?
 - Massive data challenge
 - Massive model size challenge
 - Theoretical challenge: classic ML models unable to cope

COLLABORATIVE FILTERING

Movie and product recommendation
billions model parameters



READMISSION PREDICTION

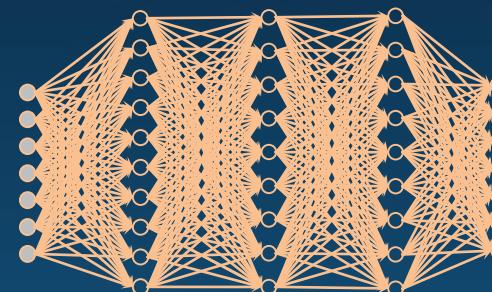
Analyzing Electronic Medical Records and
Genomics data

billions model parameters



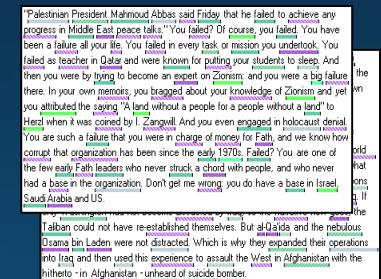
GOOGLE BRAIN

Deep learning for images
billions model parameters



BAYESIAN TOPIC MODELS

Document analysis and summarization
trillion model parameters

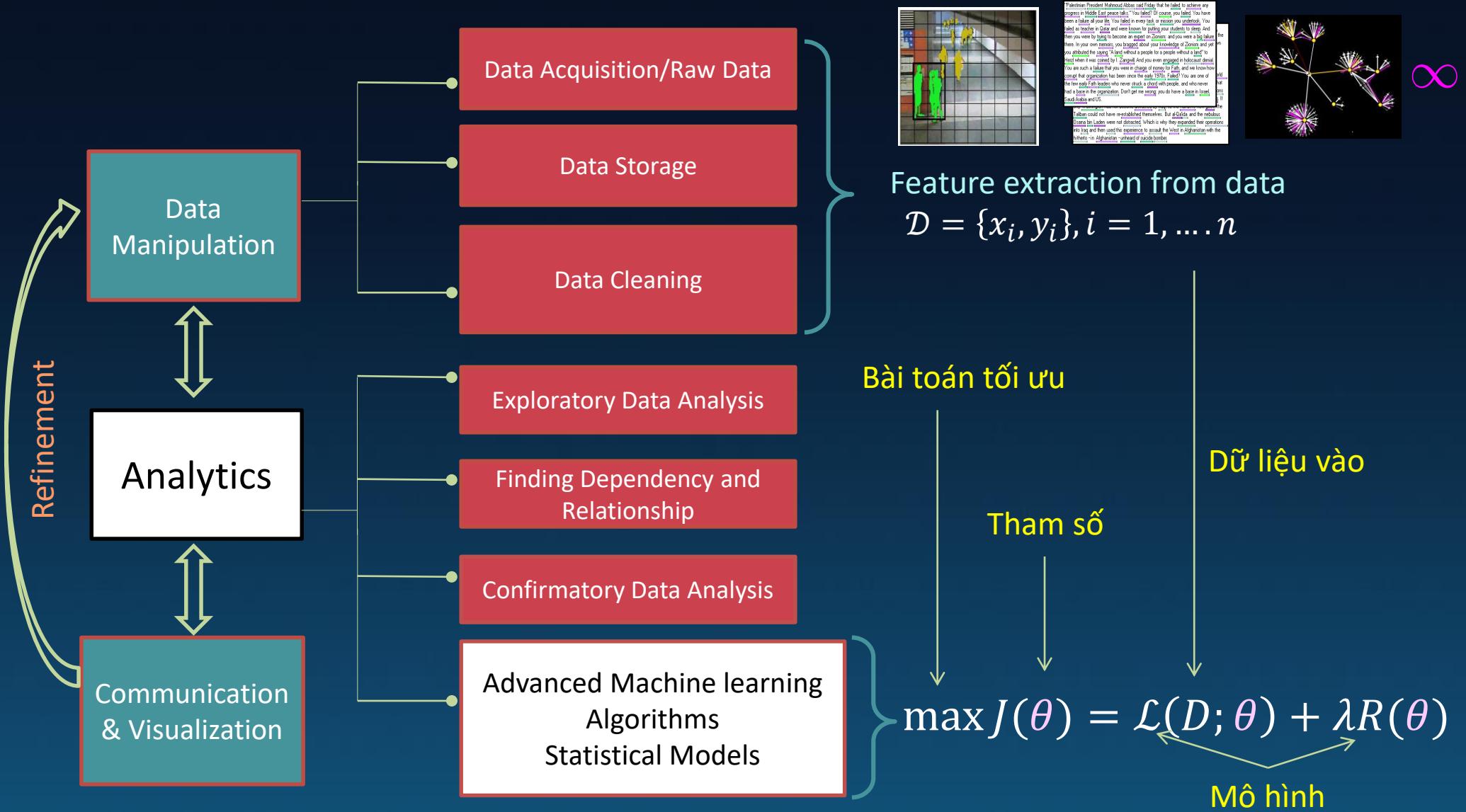


Scaling up ML and Stats models

[Xing and Qirong, Tuts KDD'15]

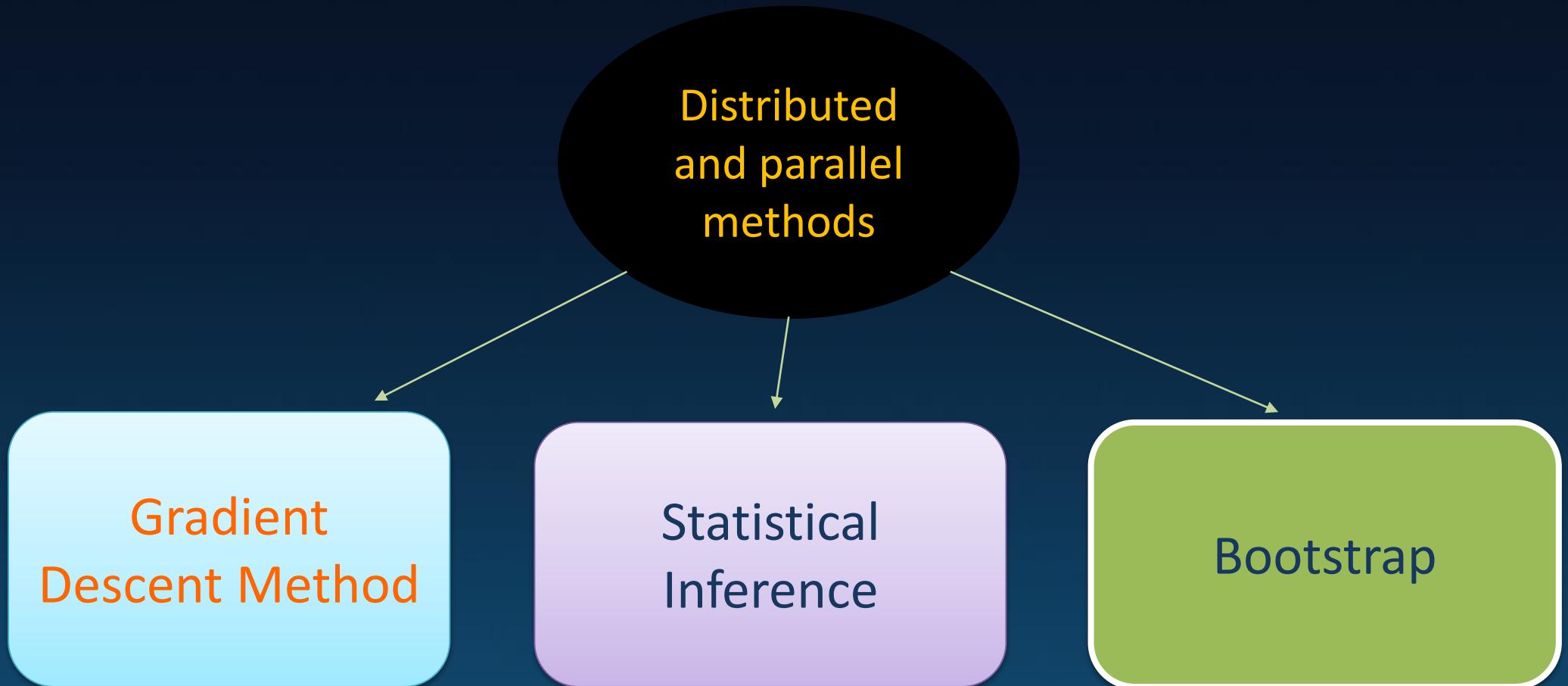
- Modern ML/Stats algorithms vs traditional program
 - Traditional programs are transaction-centric, require individual correctness at every step.
 - ML algorithms are optimization-centric and iterative convergent algorithms.
 - Error tolerance: robust against small variations/errors during intermediate steps.
 - Structural dependency: changing model parameterization and optimization strategy affects efficient parallelization and scalability.
 - Non-uniform convergence: optimal solutions might converge very differently at different number of steps.

Scale up ML/Stats models



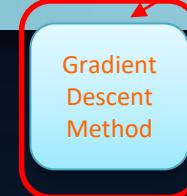
Scale up ML/Stats

- Three approaches to scale up ML/Stats models



Scale up Gradient Descent optimization

Distributed
and parallel



Statistical Inference

Bootstrap

● Distributed/Parallel Gradient Methods

- Many machine learning methods reduce to minimizing some objective function.

- Linear regression: $J(\theta) = \frac{1}{2} \sum_i [\theta^T x_i - y_i]^2$
- Logistic regression: $J(\theta) = \sum_i \log(1 + e^{-y_i \theta^T x_i})$

- Almost all objectives are optimized with **Gradient Descent (GD)** update

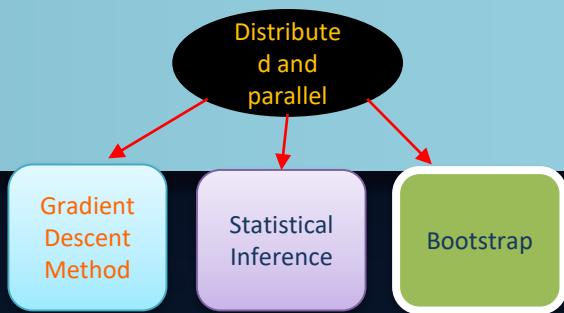
$$\theta \leftarrow \theta - \eta \sum_{i=1}^n g(\theta, x_i, y_i)$$

Iterative-convergent

- Linear regression: $g = (y_i - \theta^T x_i)x_i$

- Logistic regression: $g = y_i \left(1 + e^{-y_i \theta^T x_i}\right)^{-1} x_i$

Scale up Gradient Descent optimization



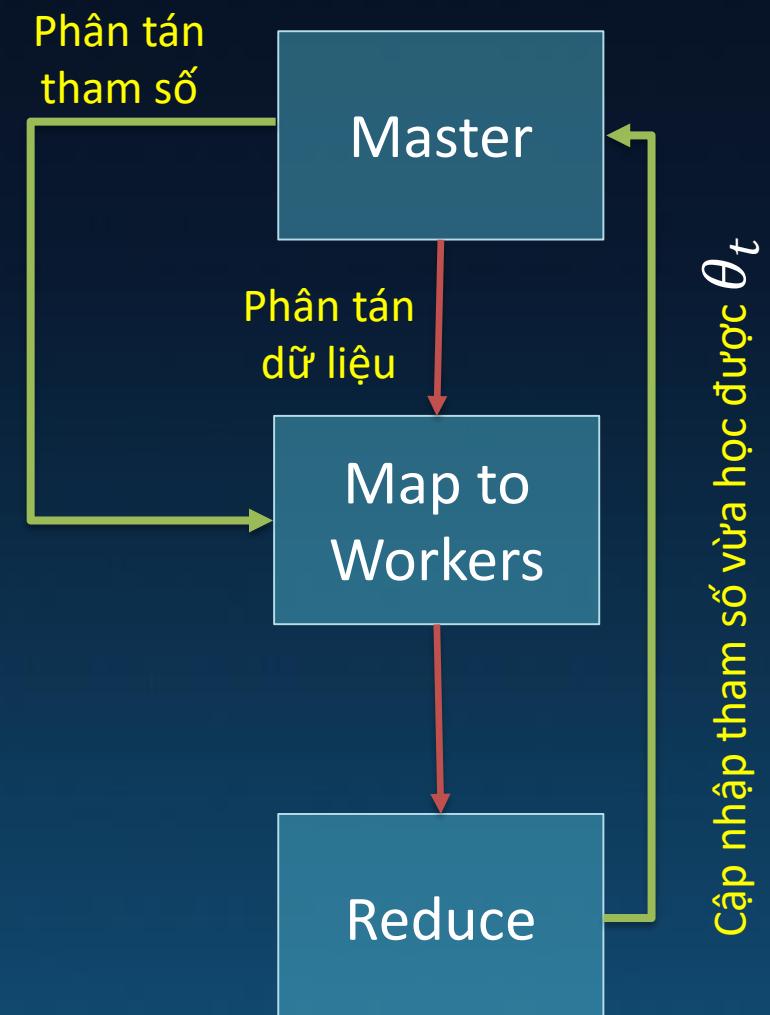
- Gradient Descent-based ML models

- Regression (linear, ridge, Lasso, ...)
- Kernel methods (SVM, SVDD, ..)
- Ensemble models (boosted gradient tree, ...)
- Almost all deep learning methods
 - Deep neural nets
 - Neural embedding: word2vec, node2vec, etc
 - Deep autoencoder: denoising autoencoder
 - Deep generative models
 - Variational Autoencoder (VAE)
 - Generative Adversarial Networks (GAN)

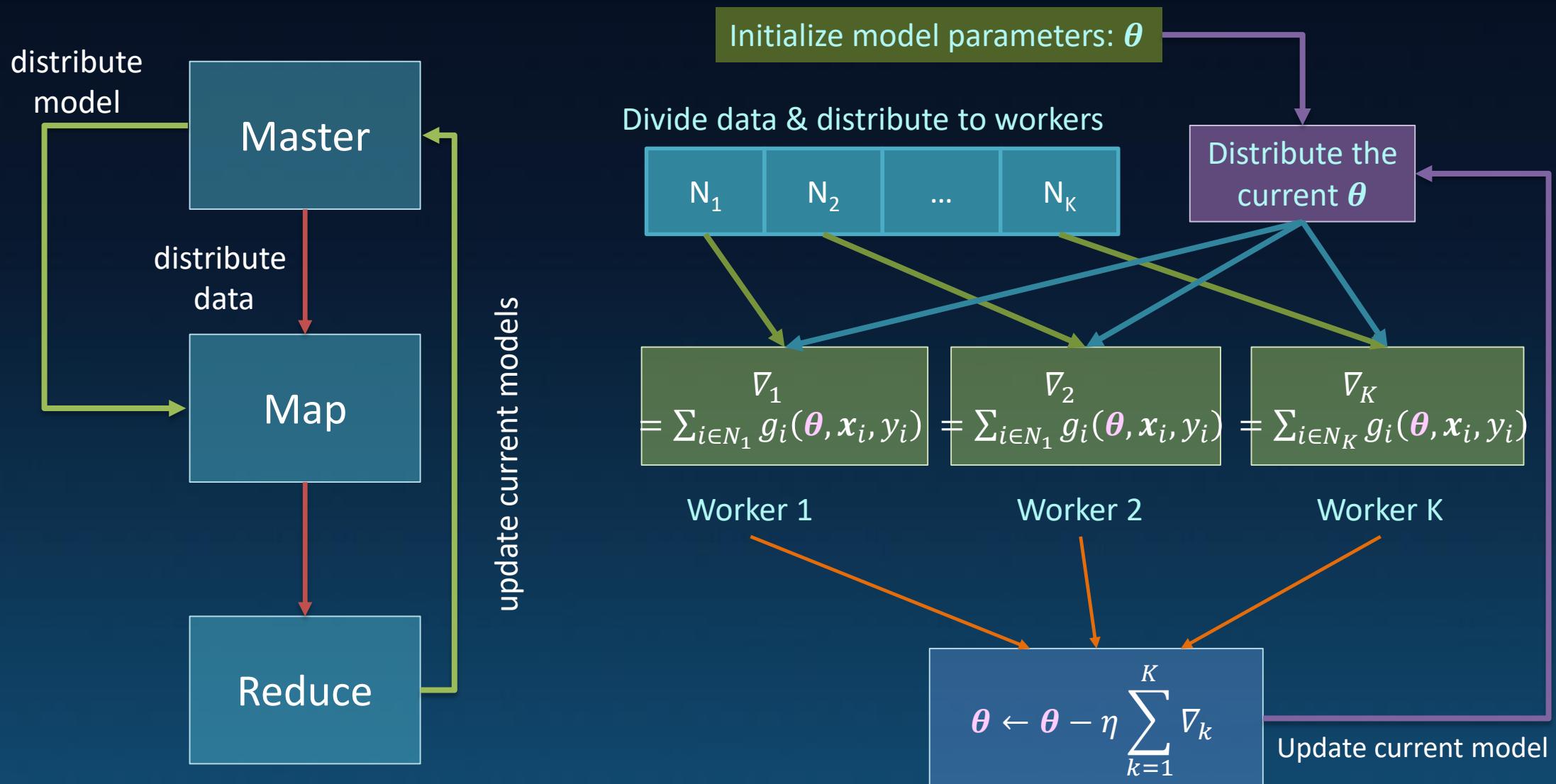
Scale up Gradient Descent optimization

- **Distributed Gradient Methods**

- While not converged
 - **Master**: distribute model parameters to **workers**.
 - **Master**: evoke map-reduce jobs
 - Many mappers compute gradient w.r.t chunks of datasets given.
 - Single reducer to combine full gradient.
 - EndWhile



Scale up Gradient Descent optimization



Scale up Gradient Descent optimization

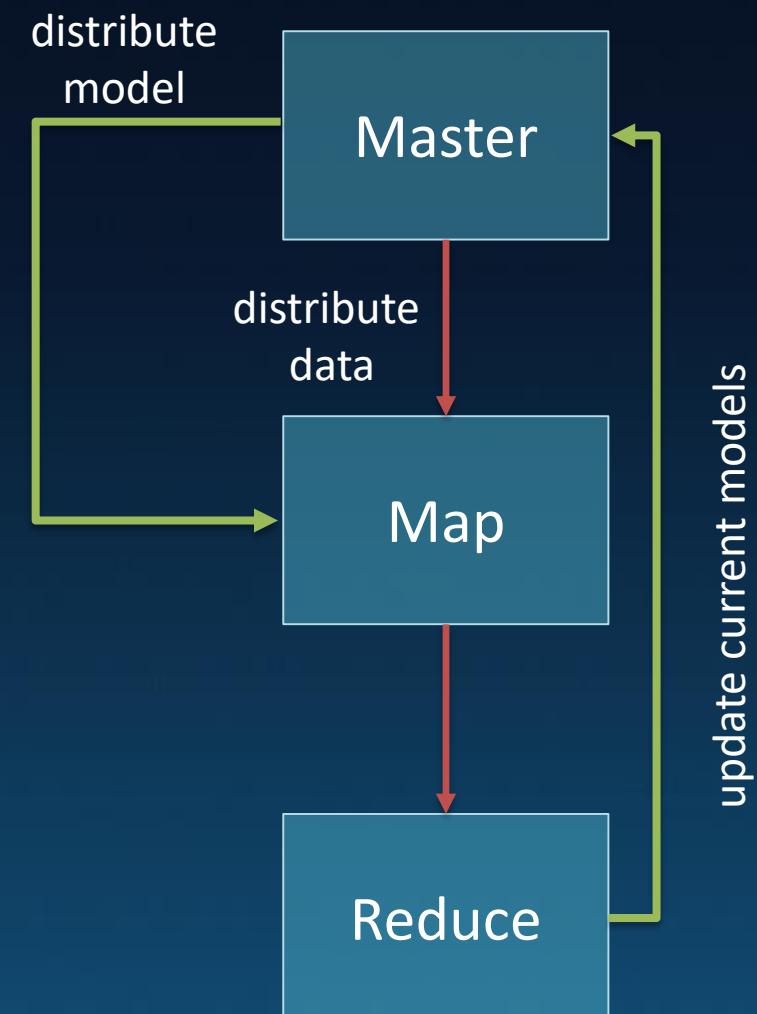
- **Distributed Gradient Methods**

- **Ưu điểm:**

- Consistency tuyệt đối
 - Hội tụ đến lời giải tốt ưu (nếu lặp đủ)
 - Tốc độ hội tụ lý thuyết nhanh (tuyến tính)

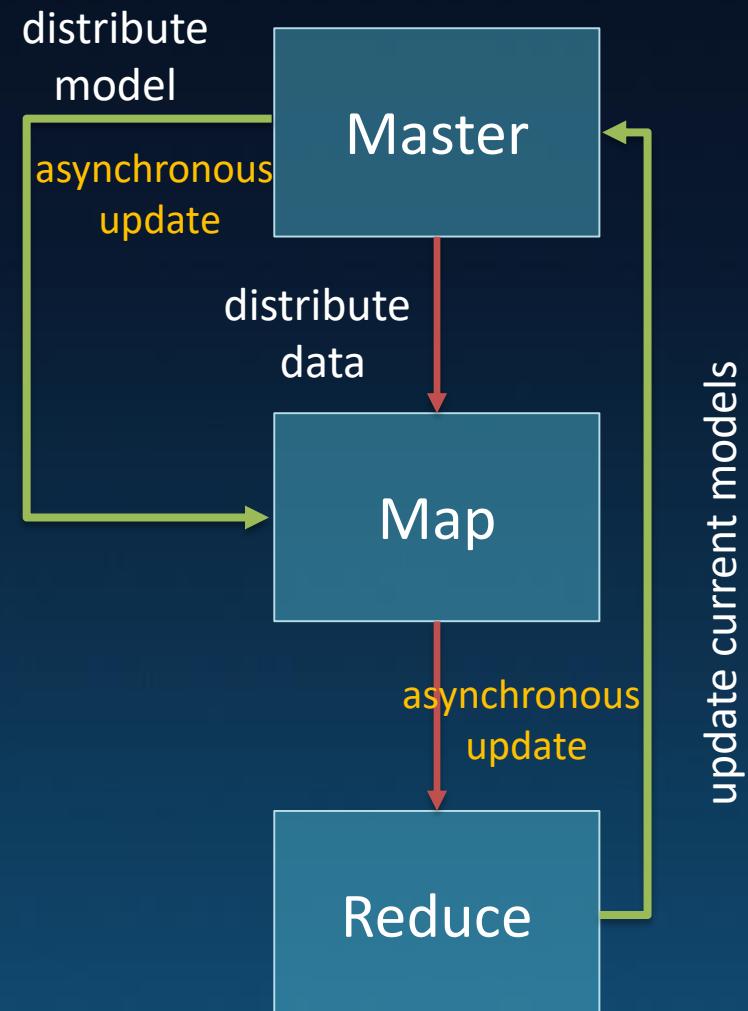
- **Nhược điểm:**

- Đòi hỏi phải có sự đồng bộ tại phía máy chủ chứa tham số, điều này khiến cho máy chủ đó phải đợi tất cả các workers tính toán xong.



Scale up Gradient Descent optimization

- Vậy làm thế nào cân bằng giữa: tính chính xác của lời giải và giảm sự đồng bộ?
- Vì sao? Giảm yêu cầu đồng bộ sẽ làm quá trình tính toán nhanh hơn.
- Chiến lược đơn giản nhất:
 - sau khi mỗi worker tính toán gradient xong thì gửi ngay lập tức về server.
 - server sẽ cập nhật tham số ngay lập tức và trả lại tham số cho worker đó.
- Chiến lược này có thể cho kết quả tốt trong một số trường hợp nhưng không đảm bảo về mặt lý thuyết.



Scale up Gradient Descent optimization

HOGWILD!: A Lock-Free Approach to Parallelizing SGD

[Recht et. al. NIPS's 11]

- Có rất nhiều bài toán tối ưu trong KHDL có thể biểu diễn lại dưới dạng như sau:

$$\min_{\mathbf{x}} f(\mathbf{x}) \triangleq \sum_{e \in E} f_e(x_e) \text{ where } f: \mathbb{R}^n \rightarrow \mathbb{R} \text{ and } E \subseteq 2^{\{1, \dots, n\}}.$$

- $|E|$ is very large and $|e|$ is small for every e

- Ví dụ như những mô hình và bài toán dưới đây:

- Sparse SVM
 - Matrix Completion
 - Graph Cuts

Scale up Gradient Descent optimization

Ví dụ: Sparse SVM

- Given a sparse training set
 - $D = \{(\mathbf{z}_1, y_1), \dots, (\mathbf{z}_N, y_N)\}$ where data item \mathbf{z}_i is sparse.
- We need to fit this data set using SVM
 - $\min_x \left(\sum_{e \in E} \max(0, 1 - y_e \mathbf{x}^T \mathbf{z}_e) + \frac{\lambda}{2} \|\mathbf{x}\|^2 \right)$ where $E = D$
- Let α_e denote the non-zero components in \mathbf{z}_e and let d_j denote the number of training examples which are non-zero in component j ($j = 1, 2, \dots, d$)

$$\min_x \left(\sum_{e \in E} \left(\max(0, 1 - y_e \mathbf{x}_{\alpha_e}^T \mathbf{z}_{\alpha_e}) + \frac{\lambda}{2} \sum_{j \in \alpha_e} \frac{x_j^2}{d_j} \right) \right)$$

Scale up Gradient Descent optimization

HOGWILD!: A Lock-Free Approach to Parallelizing SGD

Algorithm

for $t=1,2,\dots$ **do**

sample e uniformly at random from E

read current state x_e and evaluate $G_e(x_e)$

for $v \in e$ *do* $x_v = x_v - \gamma G_{ev}(x_e)$

endfor

where $G_e(x_e) = |E| \nabla f_e(x_e)$

$$\min_x f(x) \triangleq \sum_{e \in E} f_e(x_e)$$

- If the lag between when the gradient is computed and when it is used is less than a finite τ , the HOGWILD! Algorithm is proven to converge to the real optimal solution.
- This algorithm is ideal for **parallelizing in a single machine**, but less ideal for **distributing on many machines**.

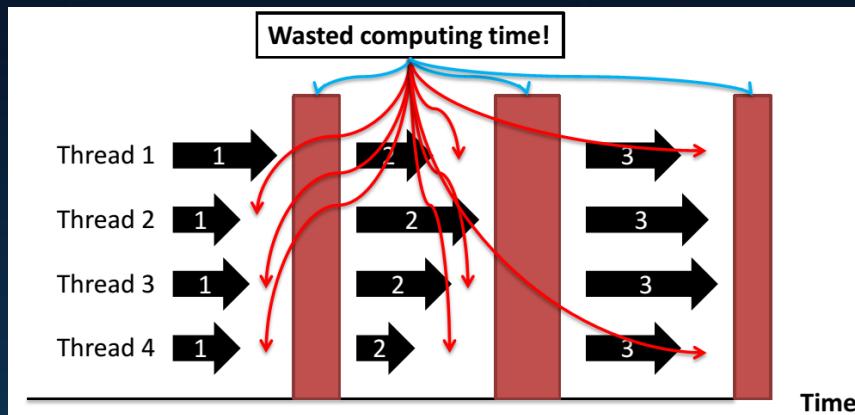
Đồng bộ hay không đồng bộ?

Synchronous vs Asynchronous

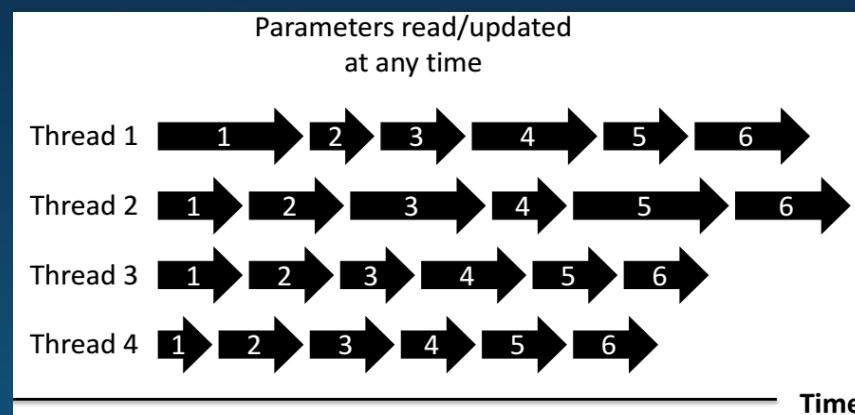
[Xing and Qirong, KDD Tuts'15]

- The fastest thread must wait for the slowest one.
- Model parameter is guaranteed to converge to optimal solution.

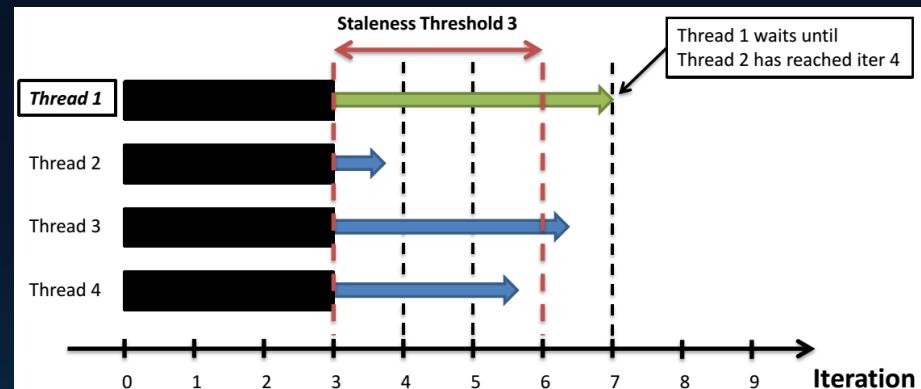
Bulk Synchronous Parallel (BSP)



Asynchronous



Stale Synchronous Parallel [Ho et. al. 2013, Dai et al 2015]

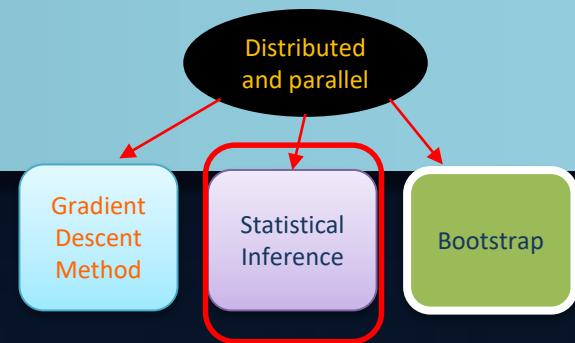
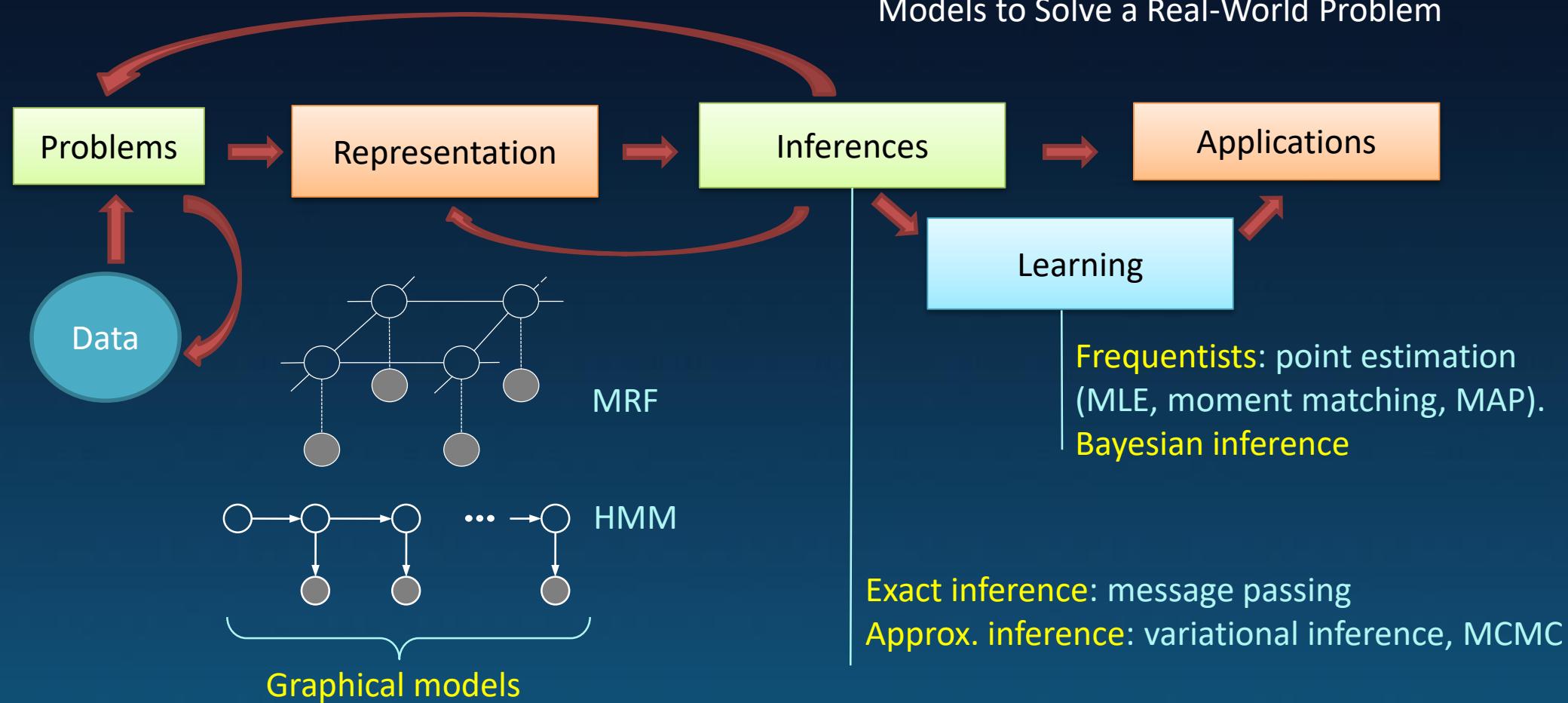


- Key idea: the time gap between fastest and slowest thread is bounded by τ .
- Theoretical guarantee is given in terms of τ .

- Model parameter is updated/read at any time whenever a thread finishes its workload.
- Model parameter might be drifted from optimal solution.

Scale up Statistical Inference

Distributed Statistical Inference



Scale up Statistical Inference

Graphical models

- Express a joint distribution over massive number of random variables.
- Conditional distributions respect graph structure $G = (V, E)$ with N vertices
- $P(x_1, \dots, x_n) = \frac{1}{Z} \prod_c \psi_c(x_c)$

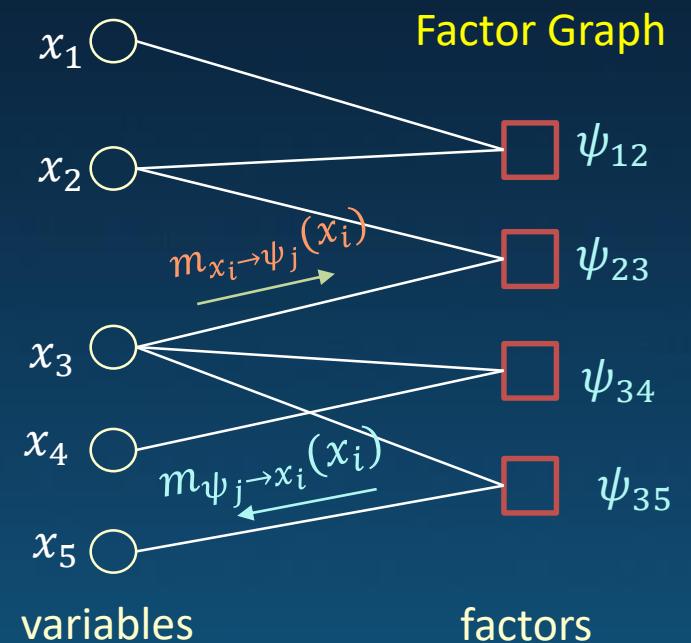
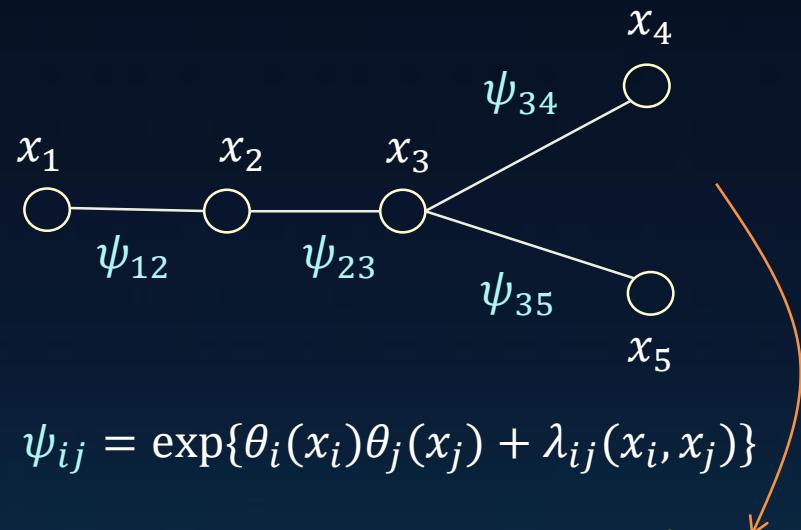
Inference = Belief Propagation

$$p(x) = \frac{1}{Z} \psi_{12} \psi_{23} \psi_{34} \psi_{35} \quad Z = \sum_{x_1, x_2, x_3, x_4, x_5} \psi_{12} \psi_{23} \psi_{34} \psi_{35}$$

$$m_{x_i \rightarrow \psi_j}(x_i) \propto \prod_{k \in \mathcal{N}_i \setminus j} m_{\psi_k \rightarrow x_i}(x_i) \quad \text{variable to factor}$$

$$m_{\psi_j \rightarrow x_i}(x_i) \propto \sum_{\text{all except } x_i} \psi_j(x_j) \prod_{k \in \mathcal{N}_j \setminus i} m_{x_k \rightarrow \psi_j}(x_k)$$

factor to variable

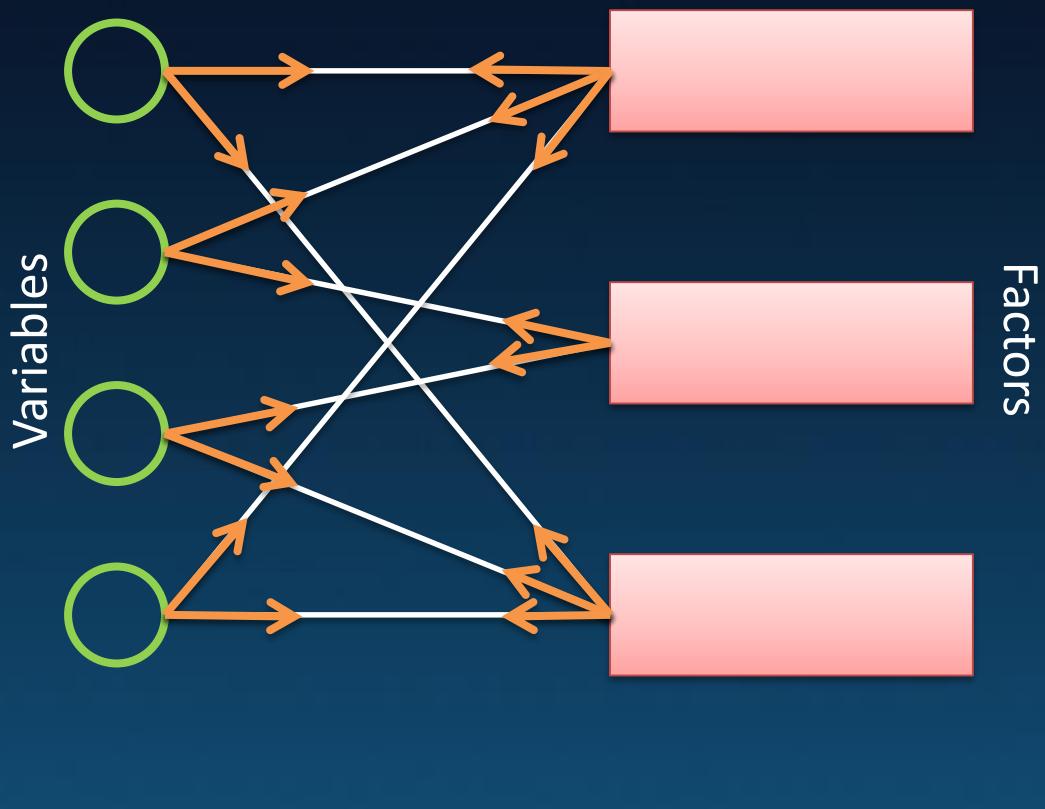


Belief Propagation

[Slides from Joseph Gonzalez et al. 09]

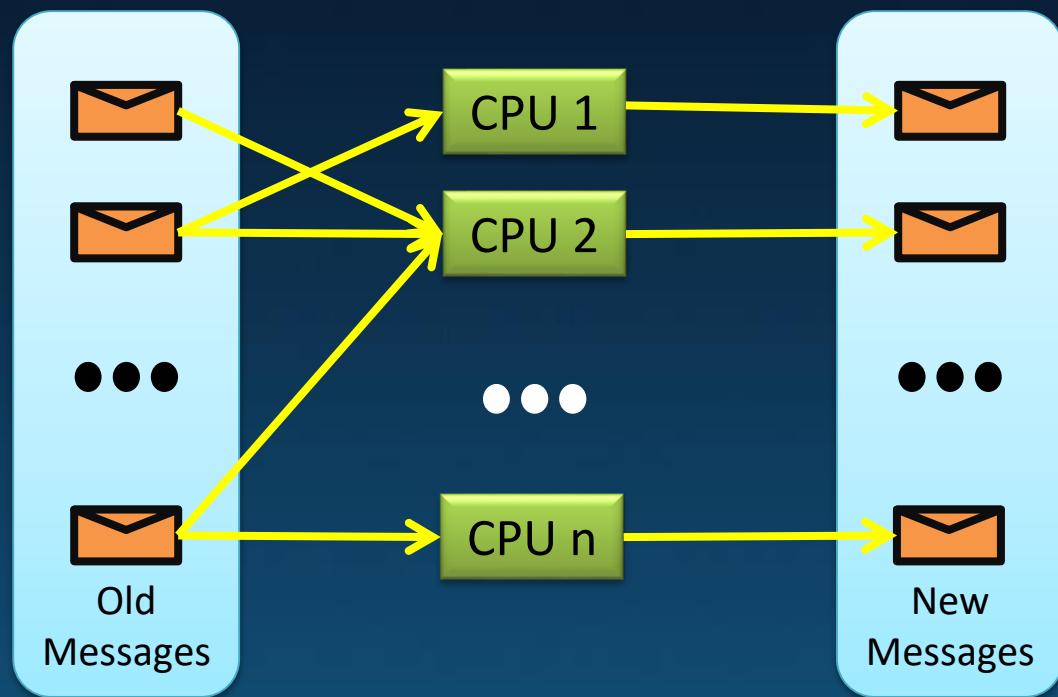
Belief Propagation (BP)

- Iterative message passing



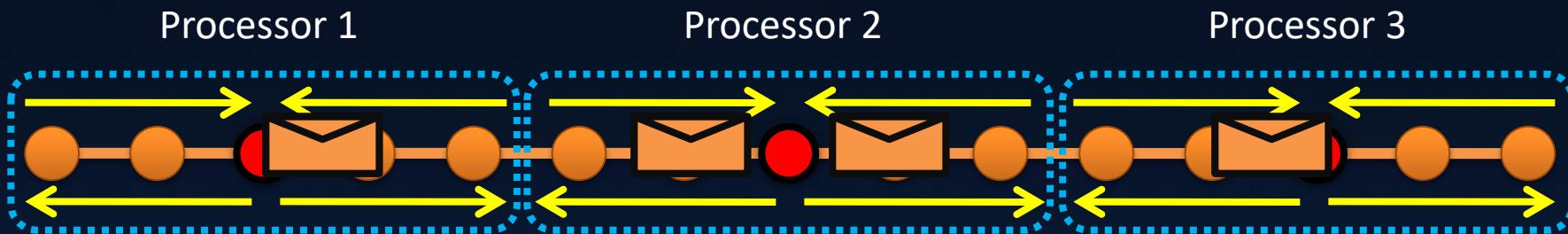
Parallel Synchronous BP

- Given all messages from previous iteration, all messages in current iteration can be computed simultaneously and independently



ChainSplash Parallel BP

[Slides from Joseph Gonzalez et al. 09]



- ChainSplash parallel BP is a parallel **approximate** inference algorithm.
- Given a chain graph, the procedure is:
 1. Divide the chain evenly among processors
 2. Loop in parallel on all processors until converged:
 - a. On each processor, run sequential Forward – Backward
 - b. On each processor, exchange messages with its 2 neighbours

Theorem: Using p processors this algorithm achieves a τ_ϵ approximation in time:

$$O\left(\frac{n}{p} + \tau_\epsilon\right)$$

The Splash Belief Propagation algorithm

Splash Belief
Propagation



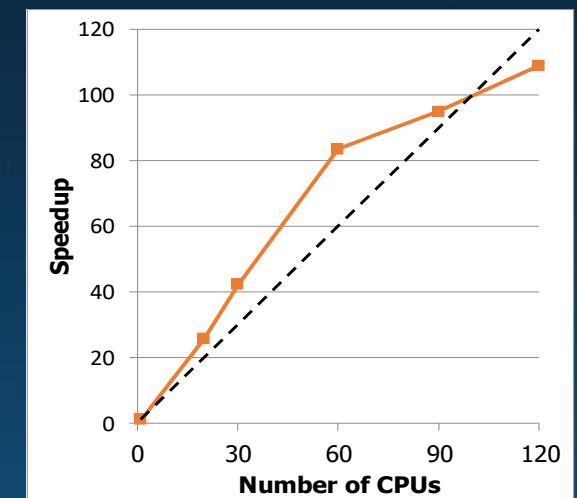
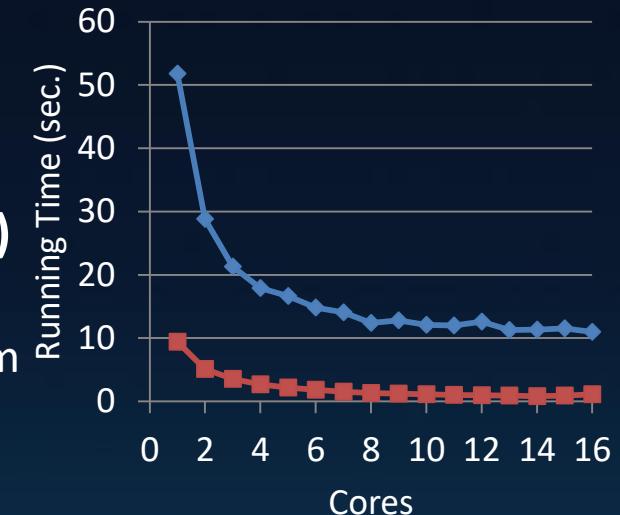
Splash
operation



dynamic
Splash
scheduling

Splash BP (orange)
speedup performance
on distributed system with
a large graph

Synchronous BP (blue)
vs Splash BP (red)
on shared memory system
with a small graph

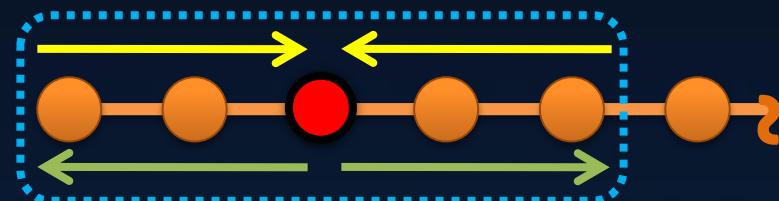


[Nguồn: Joseph Gonzalez et al., 2009]

The Splash Operation

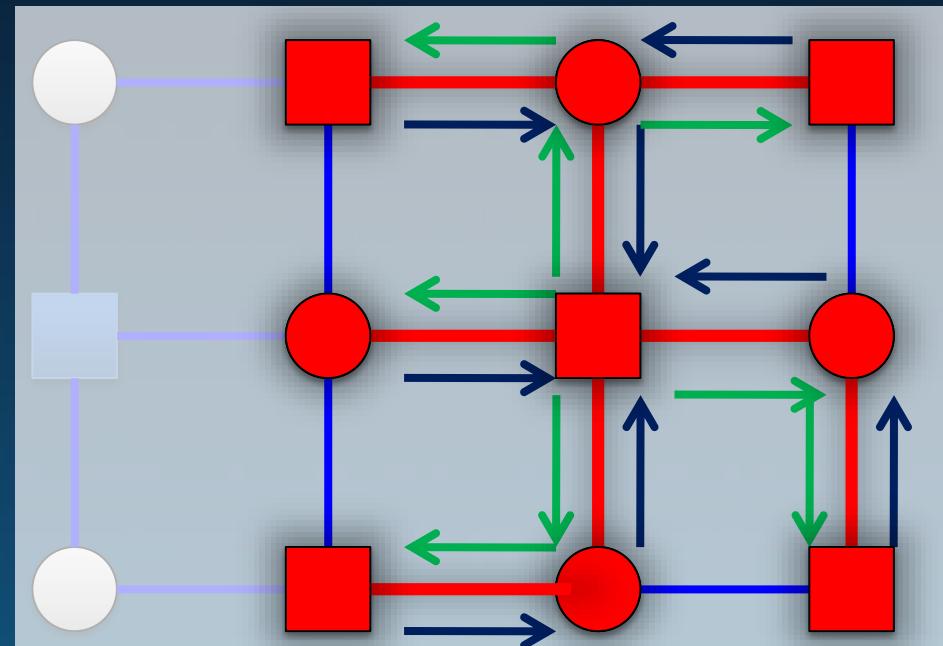
[Slides from Joseph Gonzalez et al. 09]

- Generalize the optimal chain algorithm:



to arbitrary cyclic graphs:

- Grow a BFS Spanning tree with fixed size.
- Forward Pass computing all messages at each vertex.
- Backward Pass computing all messages at each vertex.

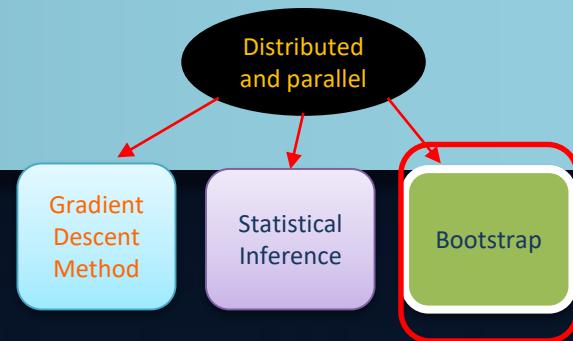


Scale up Statistical Inference

- Một số cách tiếp cận khác
 - MCMC inference for graphical models
 - Naïve approach: take multiple chains in parallel, take average/consensus between chains.
 - Embarrassingly Parallel (but correct) MCMC [Neiswanger et al, 14]
 - Parallel Gibbs sampling
 - Approximated distributed LDA [Newman et al, 2009]
 - Asynchronous version [Ahmed et al. 2012, Dai et al. 2015]
 - see Eric Xing's KDD 2015 tutorials for more references.
 - Variational inference for graphical models
 - Stochastic Variational Inference (SVI) [Hoffman et al., JMLR13, Huynh et al, UAI16]
 - Deterministic Wasserstein-based Optimization [Nguyen, Annals of Stats, 13]

Scale up Model Evaluation

Bootstrap for massive data [Kleiner, Jordan 12]



- It is an important theoretical aspect to assess the quality of inference.

Data $X_1, X_2, \dots, X_n \sim f_\theta$ —————> Estimate $\hat{\theta}$ —————> Quality $|\theta - \hat{\theta}|$?

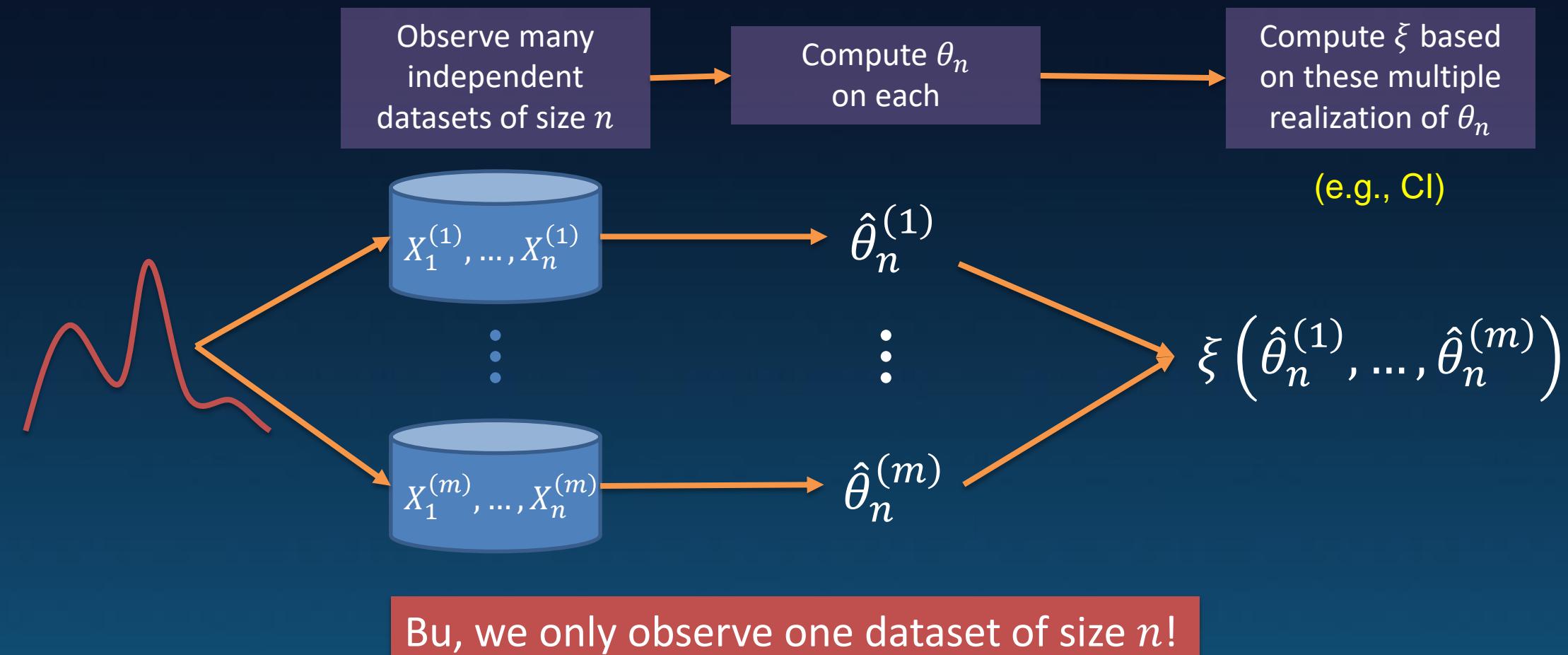
- Bootstrap: generic framework for computing confidence interval
 - it's useful in ML/Stats to assess the quality of the estimates (e.g., model parameters) in terms of confidence interval, error bars.
 - Model selection, smoothing, bias correction, ...
 - Provide more useful information than a simple point estimate.
- Can it be used on big data?

Scale up Model Evaluation

Bootstrap for massive data

[Kleiner, Jordan 12]

- Trường hợp lý tưởng:



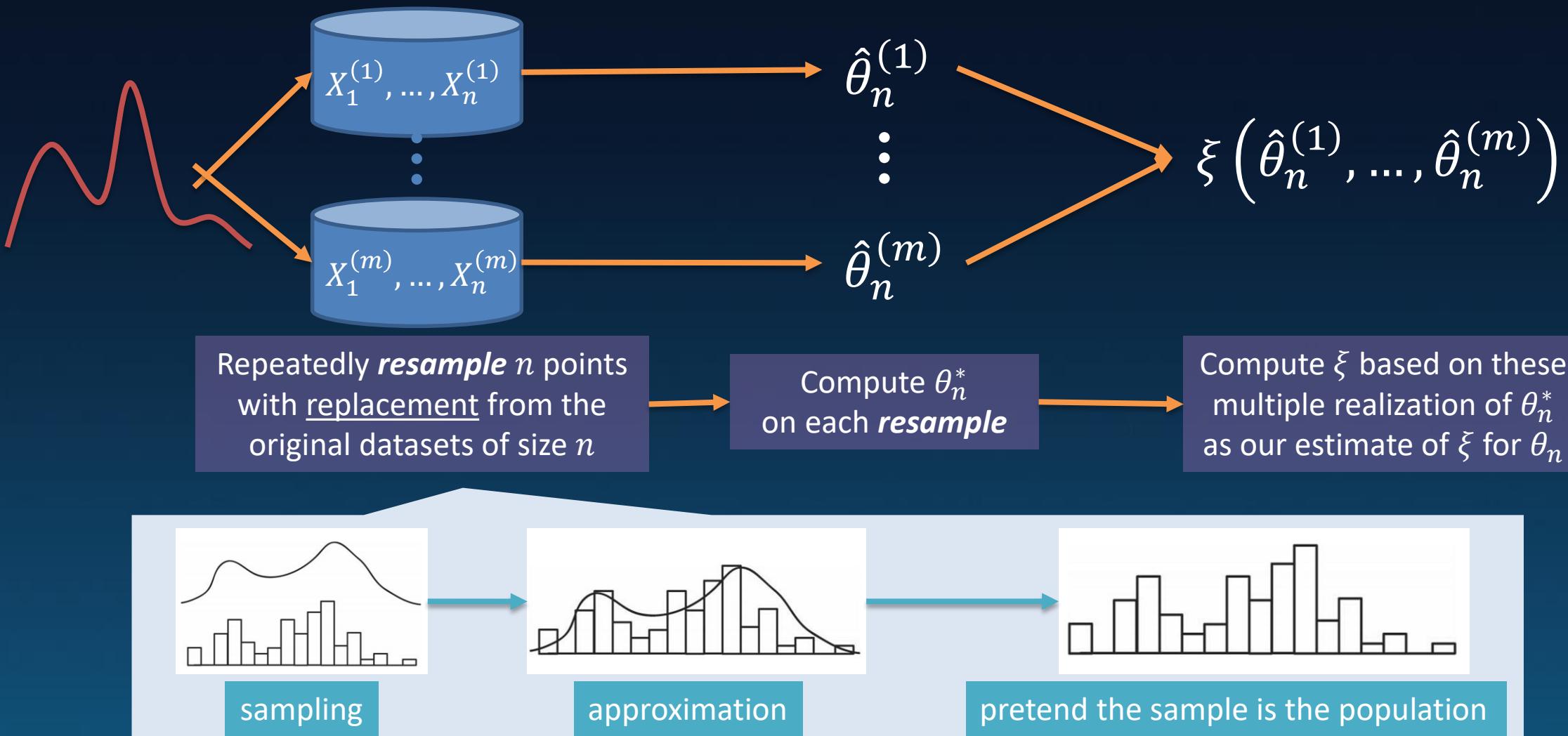
Scale up Model Evaluation

Bootstrap for massive data

[Kleiner, Jordan 12]

Bootstrap (Efron, 1979)

- Use the observed data to simulate datasets of size n



Scale up Model Evaluation

Bootstrap for massive data

- Computational bottleneck with bootstrap

[Slide from Jordan, ICML14]

- Seemingly a wonderful match to modern parallel and distributed computing platforms
 - But the expected number of distinct points in a bootstrap resample is $\sim 0.632n$

○ E.g., if original dataset has size 1TB, then expect resample to have size $\sim 632GB$

- Can't feasibly send resampled datasets of this size to distributed servers.
 - Even if one could, can't compute the estimate locally on datasets this large.

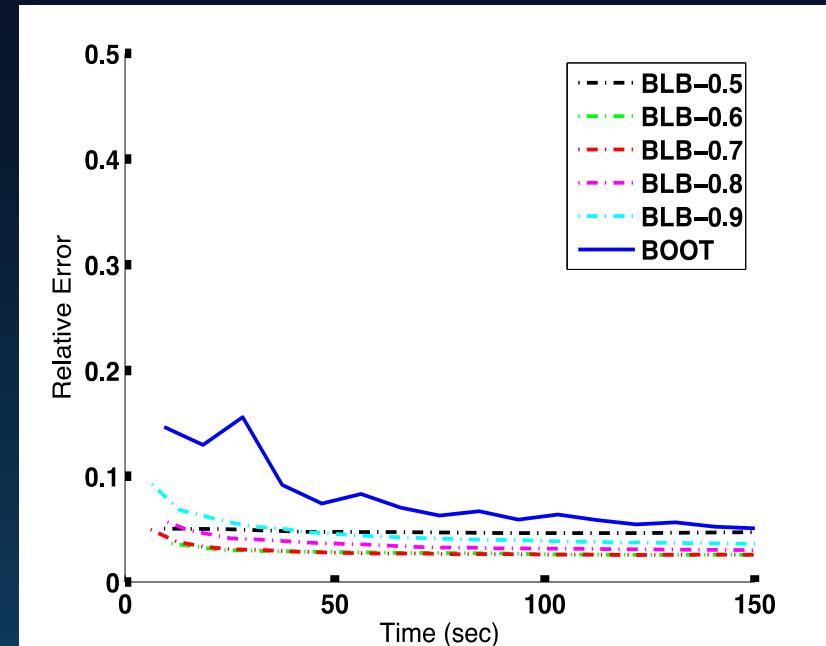
Scale up Model Evaluation

Bootstrap for massive data

[Kleiner, Jordan 12]

• The Bag of Little Bootstrap (BLB)

- The subsample contains only b points, and so the resulting empirical distribution has its support on b points
- But we can (and should!) resample it with replacement n times, not b times
- Doing this repeatedly for a given subsample gives bootstrap confidence intervals on the right scale---no analytical rescaling is necessary!
- Now do this (in parallel) for multiple subsamples and combine the results (e.g., by averaging)



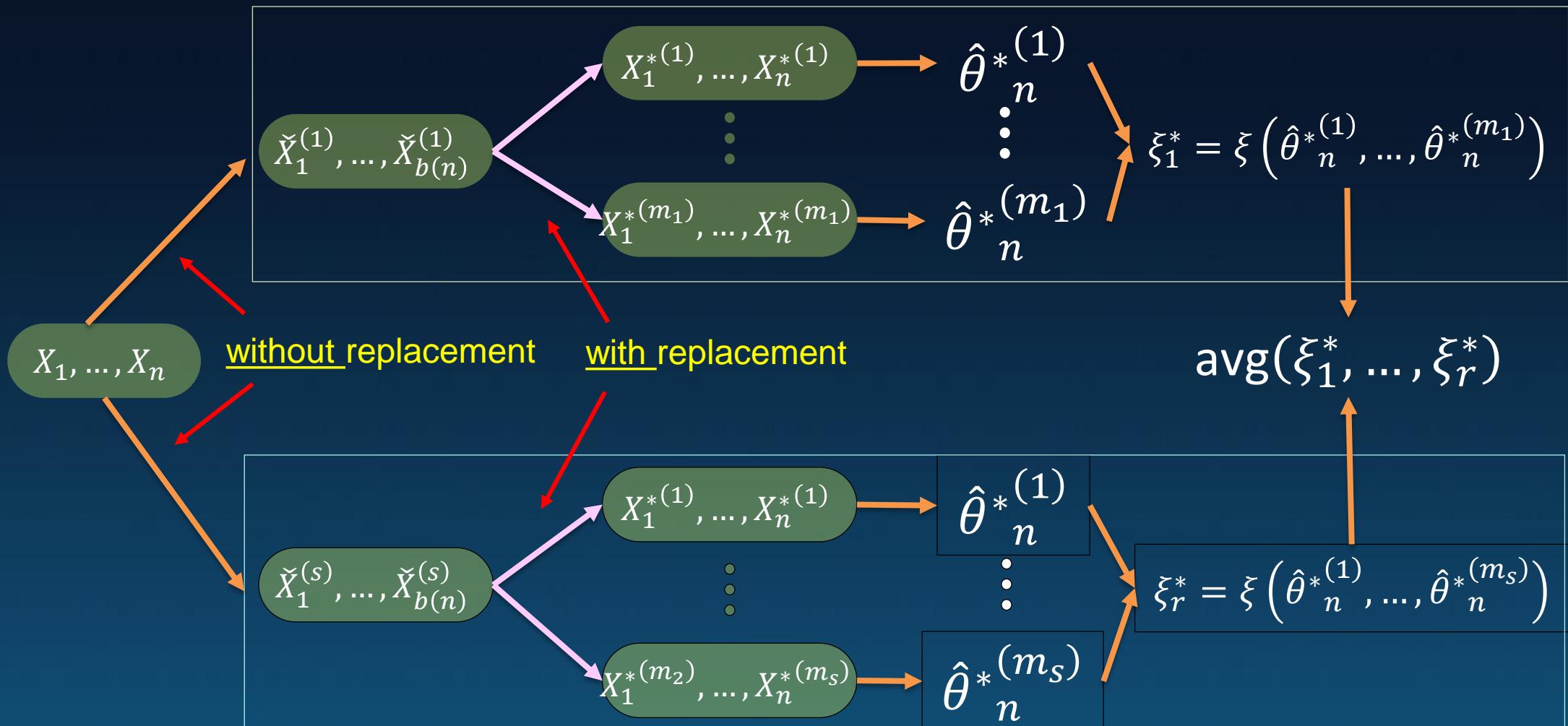
[Slides from Jordan, ICML14]

Scale up Model Evaluation

Bootstrap for massive data

[Kleiner, Jordan 12]

- The Bag of Little Bootstrap (BLB)



Scale up Model Evaluation

Bootstrap for massive data

[Kleiner, Jordan 12]

- BLB is asymptotically consistent and higher-order correct (like the bootstrap), under essentially the same conditions that have been used in prior analysis of the bootstrap.

Theorem (asymptotic consistency): Under standard assumptions (particularly that θ is Hadamard differentiable and ξ is continuous), the output of BLB converges to the population value of ξ as n, b approach ∞ , then:

$$\left| m_1^{-1} \sum_{j=1}^{m_1} \xi \left(Q_n \left(\mathbb{P}_{n,b}^{(j)} \right) \right) - \xi(Q_n(P)) \right| = O_P \left(\frac{\sqrt{\text{Var}(\hat{p}_k^{(j)} - p_k | \mathbb{P}_n)}}{\sqrt{nm_1}} \right) + O_P \left(\frac{1}{n} \right) + O \left(\frac{1}{b\sqrt{n}} \right)$$

Therefore, taking $m_1 = \Omega(n \text{Var}(\hat{p}_k^{(j)} - p_k | \mathbb{P}_n))$ and $b = \Omega(\sqrt{n})$ yields

$$\left| m_1^{-1} \sum_{j=1}^{m_1} \xi \left(Q_n \left(\mathbb{P}_{n,b}^{(j)} \right) \right) - \xi(Q_n(P)) \right| = o_P \left(\frac{1}{n} \right),$$

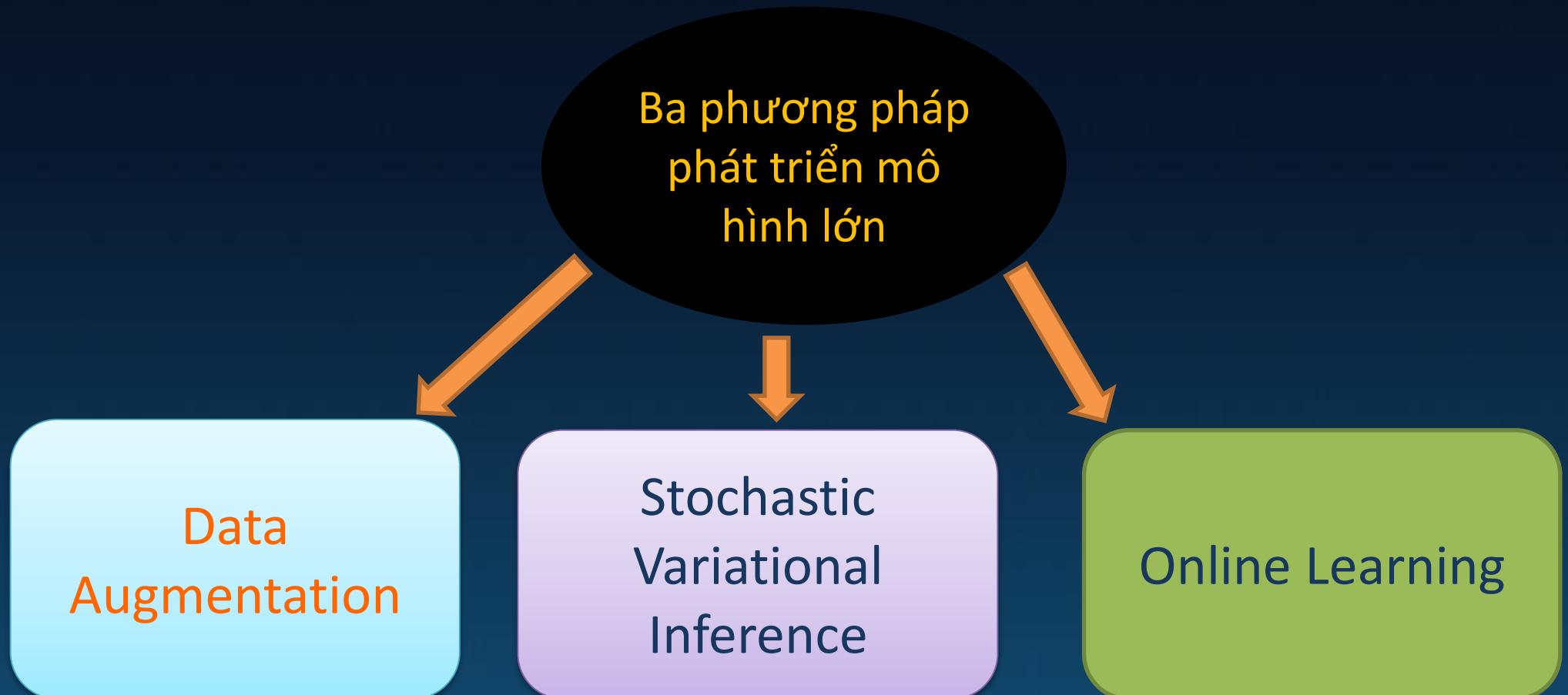
in which case BLB enjoys the same level of higher-order correctness as the bootstrap.

Ba cách tiếp cận để phát triển mô hình lớn

Three approaches to develop
new scalable models

- Mô hình lớn cho dữ liệu lớn
 - Three approaches to scale up big models
 - Scale up Gradient Descent Methods
 - Scale up Probabilistic Inference
 - Scale up Model Evaluation (bootstrap)
- Ba cách tiếp cận để phát triển mô hình lớn
 - Data augmentation
 - Stochastic Variational Inference for Graphical Models
 - Stochastic Gradient Descent and Online Learning
- Công nghệ nào cho mô hình lớn?
 - Hadoop, Spark, TensorFlow

Three Approaches to Develop Scalable Models



Data Augmentation

Phương pháp bổ sung tham số độc lập

Ba phương pháp
phát triển mô hình lớn

Data
Augmentation

Stochastic
Variational
Inference

Online
Learning

- Các bước chính:

- 1) Formulate OP as sampling problem
- 2) Aim to sample the joint distribution via Gibbs sampling
- 3) Introduce auxiliary variables to turn intractable conditional distributions into tractable forms

- A useful result [Polson JMLR'13]: from Hinge-loss to distribution

- $\exp(-|ab|) = \int_0^\infty \frac{a}{\sqrt{2\pi\lambda}} \exp\left\{-\frac{1}{2}(a^2\lambda + b^2\lambda^{-1})\right\} d\lambda$
- Hence, $\exp\{-2\max(1 - y_i \boldsymbol{\theta}^T \mathbf{x}_i)\} = \int_0^\infty \frac{1}{\sqrt{2\pi\lambda_i}} \exp\left\{-\frac{(1+\lambda_i-y_i \boldsymbol{\theta}^T \mathbf{x}_i)^2}{2\lambda_i}\right\} d\lambda_i$
- Therefore, $p(y_i, \lambda_i | \boldsymbol{\theta}, \mathbf{x}_i) = \frac{1}{\sqrt{2\pi\lambda_i}} \exp\left\{-\frac{(1+\lambda_i-y_i \boldsymbol{\theta}^T \mathbf{x}_i)^2}{2\lambda_i}\right\}$ We now can sample!

Data Augmentation

Giải bài toán SVM với hàng trăm triệu điểm dữ liệu

[Tu et al, ICPR16]

• Cách tiếp cận

Chuyển bài toán tối ưu sang suy diễn Bayes (MAP)

Đại lượng khó xử lý
(intractable)

Dùng kết quả của Polson để thêm biến phụ độc lập (auxiliary variables)

Bây giờ ta đã có thể tính song song

SVM giải bài toán tối ưu

$$\min_{\boldsymbol{\theta}} \left(\frac{\alpha}{2} \|\boldsymbol{\theta}\|_2^2 + \frac{1}{N} \sum_{n=1}^N \max(0, 1 - y_n \boldsymbol{\theta}^T \mathbf{x}_n) \right)$$

The conditional distribution for $\boldsymbol{\theta}$ is

$$p(\boldsymbol{\theta} | X, \mathbf{y}, \alpha) \propto p(\mathbf{y} | X, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \alpha)$$

where

- $p(\mathbf{y} | X, \boldsymbol{\theta}) \propto \frac{1}{N} \sum_n \exp(-\max(0, 1 - y_n \boldsymbol{\theta}^T \mathbf{x}_n))$
- $p(\boldsymbol{\theta} | \alpha) = \mathcal{N}(0, \alpha^{-1})$

Introduce auxiliary variables $\lambda_{1:N}$ such that

$$\int_{\lambda_n} p(y_n, \lambda_n | \mathbf{x}_n, \boldsymbol{\theta}) = p(y_n | \mathbf{x}_n, \boldsymbol{\theta})$$

- Sample $\boldsymbol{\theta}$:

$$p(\boldsymbol{\theta} | X, \mathbf{y}, \boldsymbol{\tau}, \alpha) \propto \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- Sample λ_n :

$$p(\lambda_n | \boldsymbol{\theta}, X, \mathbf{y}, \boldsymbol{\tau}_{-n}, \lambda) \sim \mathcal{IG}(|1 - y_n \boldsymbol{\theta}^T \mathbf{x}_n|^{-1}, 1)$$

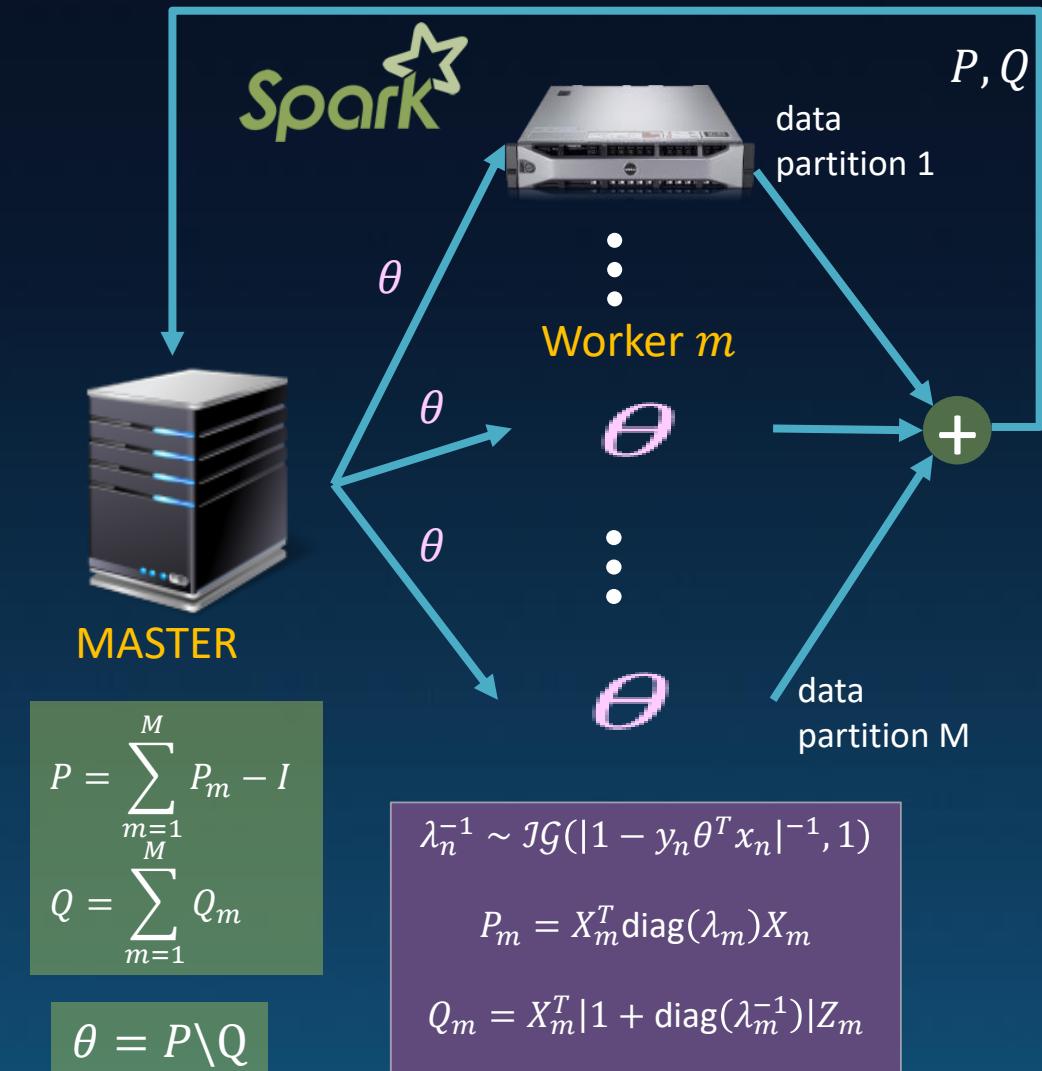
Data Augmentation

Giải bài toán SVM với hàng trăm triệu điểm dữ liệu

[Tu et al, ICPR16]

Distributed Data Augmented Support Vector Machine

- Sampling each λ_n is independent
- Sample θ :
 - We have $p(\theta | X, y, \tau, \lambda) \propto \mathcal{N}(\mu, \Sigma)$.
 - Then $\theta_{MAP} = \mu$ where
$$\mu = \Sigma[X^T(1 + \lambda^{-1})]$$
$$\Sigma^{-1} = X^T \text{diag}(\lambda^{-1})X$$
 - Let $P = X^T \text{diag}(\lambda)X + I$ and $Q = X^T Z$ where $Z = \{-1, 1\}^{NK \times 2}$ and $z_{nk}=1$ if $y_n = (k+1)/2$, then we have $\theta = P \setminus Q$ (i.e, $P\theta_k = Q_k$)
 - Can be parallelized in a distributed system in a similar way to parallelize matrix multiplication problem



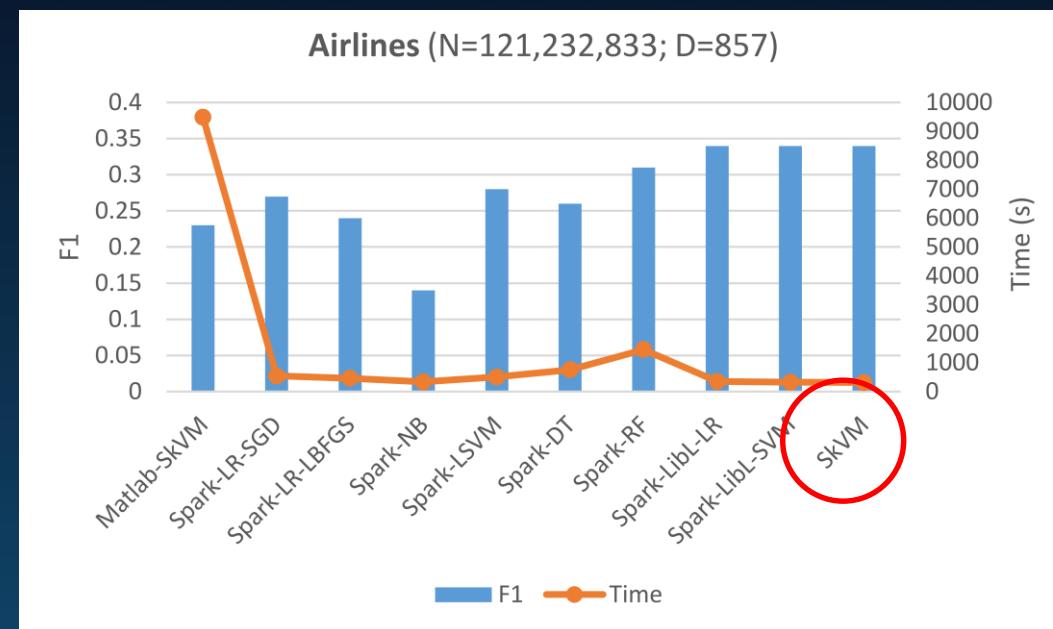
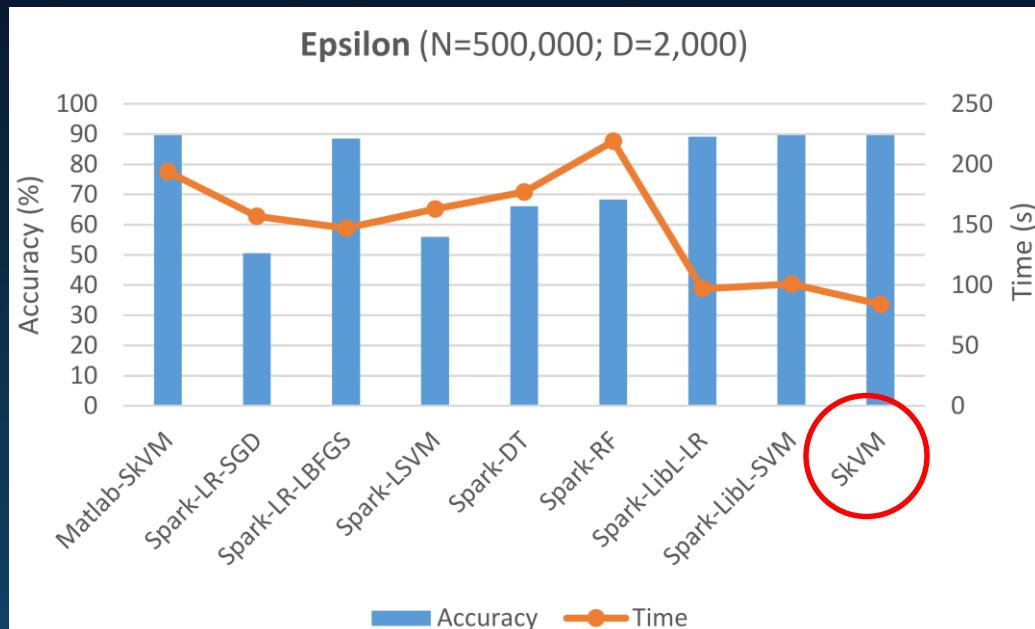
Data Augmentation

Giải bài toán SVM với hàng trăm triệu điểm dữ liệu

[Tu et al, ICPR16]

- Tập dữ liệu:

- Epsilon: # data points = 500,000 and dimension = 2,000
- Airline: # data points = 121,232,833 and dimension = 857



- For scaling up logistic regression, see our ICDM'16 paper.

Variational Inference

Suy diễn biến phân

Ba phương pháp
phát triển mô hình lớn

Data
Augmentation

Stochastic
Variational
Inference

Online
Learning

- Học một mô hình θ từ dữ liệu x bằng phương pháp suy diễn Bayes cũng chính là tìm phân bố posterior $p(\theta | x)$

- Variational inference

- Approximate $p(\theta | x)$ with a $q(\theta)$ where $q(\theta)$ is tractable
 - How to evaluate the closeness: Kullback-Leibler (KL) divergence
$$KL(q, p) = \int q(\theta) \ln \frac{q(\theta)}{p(\theta | x)} d\theta$$
 - Optimize KL divergence from $q(\theta)$ to $p(\theta | D)$ to is to maximize the Evidenced Lower Bound (ELBO)

$$\mathcal{F}(q) = \mathbb{E}_q [\ln p(x | \theta)] - KL(q(\theta), p(\theta | x))$$

Variational Inference

Suy diễn biến phân

• Mô hình chủ đề (topic models)

$$\mathbb{P} \left(\mathbf{x} = \text{[Text Content]} \right) = \text{[Topic Model Formula with Theta and Topic Distribution Matrix]}$$

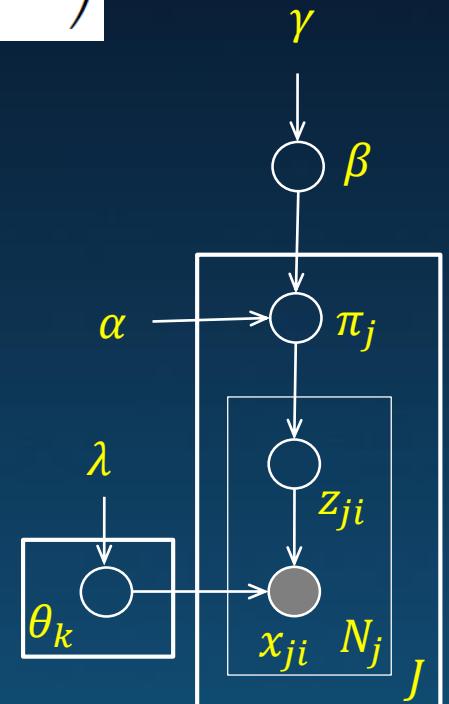
The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FELLOWSHIP	PROGRAM	WOMEN	STUDENTS
SHOW	BUDGET	PEOPLE	SCHOOL
MUSIC	BILLION	CHILD	EDUCATION
MOVIE	FEDERAL	YEARS	TEACHERS
PLAY	YEAR	FAMILIES	HIGH
MUSICAL	SPENDING	WORK	PUBLIC
BEST	ACTOR	PARENTS	TEACHER
ACTOR	NATIONAL	SANCTUARY	BEST
FIRST	STATE	FAMILY	MANGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

• Biến ẩn (latent variables): $\Phi = \{z, \pi, \theta\}$

- Posterior: $p(\Phi | x)$, variational distribution $q(\Phi)$
- Bài toán suy diễn biến phân trở thành bài toán tìm cực đại:

$$-KL(q, p) = \int_{\Phi} q \ln \frac{p}{q} d\Phi = \mathbb{E}_q[\ln p(\Phi, x)] - \mathbb{H}(q) + C$$



Variational Inference

Suy diễn biến phân cho mô hình chủ đề

- Giả thuyết mean-field

- $q(\Phi) = q(\beta \mid \lambda^\beta) \prod_{ji} q(z_{ji} \mid \mu_{ji}^z) \prod_j q(\pi_j \mid \lambda_j^\pi) \prod_k q(\theta_k \mid \lambda_k^\theta)$

- Giải thuật suy diễn biến phân cho LDA

Loop until convergence

- For $j = 1: J$

- For $i = 1: N_j$

- Update $q(z_{ji}) \propto \exp\left(\mathbb{E}_{q-z_{ji}}[\ln p(\Phi, x_{ji})]\right)$

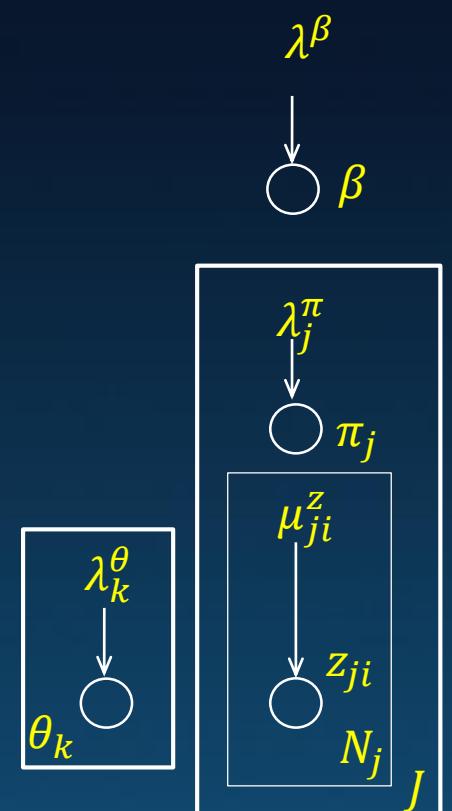
- Update $q(\pi_j) \propto \exp\left(\mathbb{E}_{q-\pi_j}[\ln p(\Phi, x_j)]\right)$

- Update $q(\theta_k) \propto \exp\left(\mathbb{E}_{q-\theta_k}[\ln p(\Phi, x)]\right)$

- Update $q(\beta) \propto \exp(\mathbb{E}_{-\beta}[\ln p(\Phi, x)])$

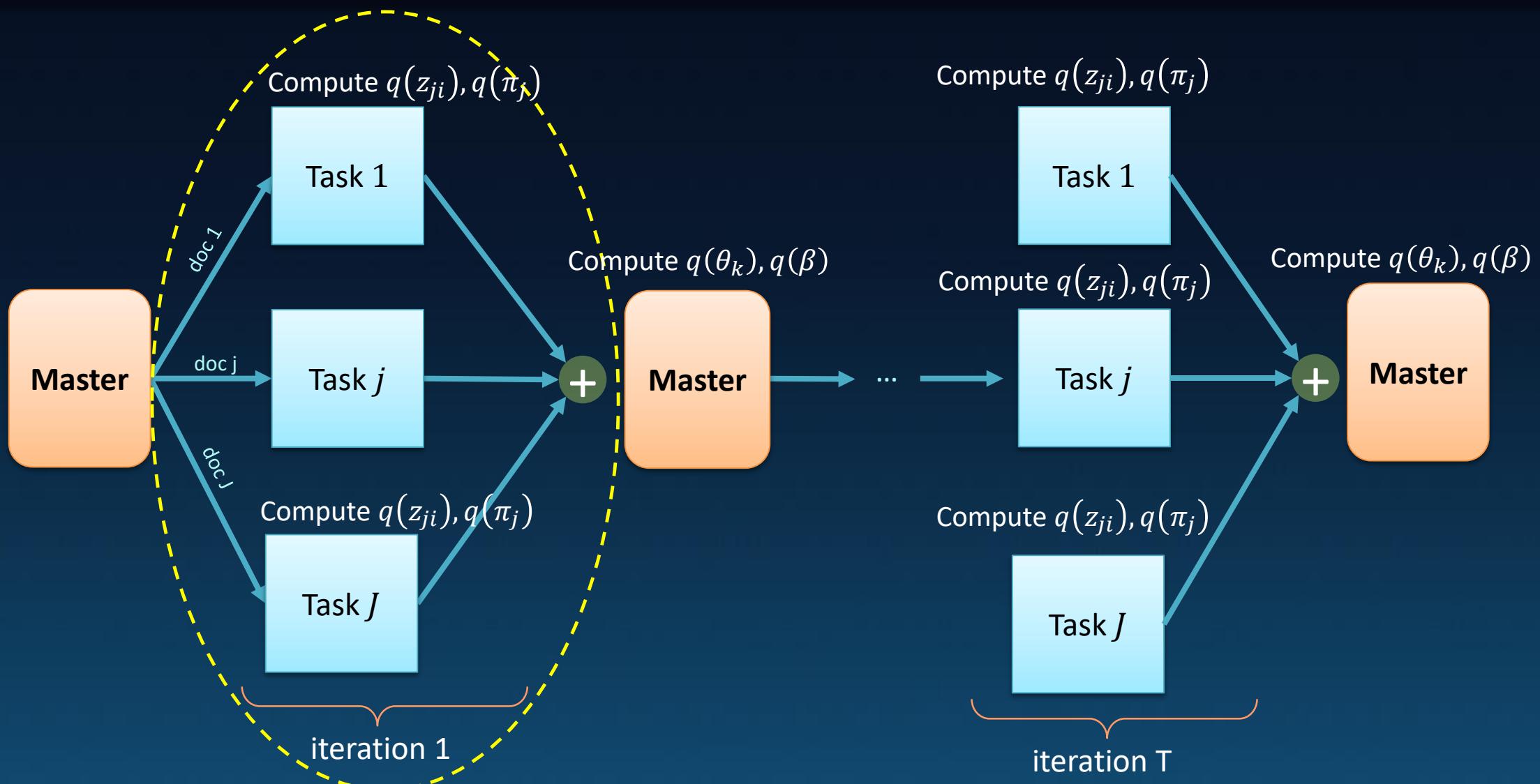
local variables

global variables



Distributed Variational Inference

Xử lý phân tán cho suy diễn biến phân cho mô hình chủ đề



Each iteration involves the computation on the entire dataset of J documents.

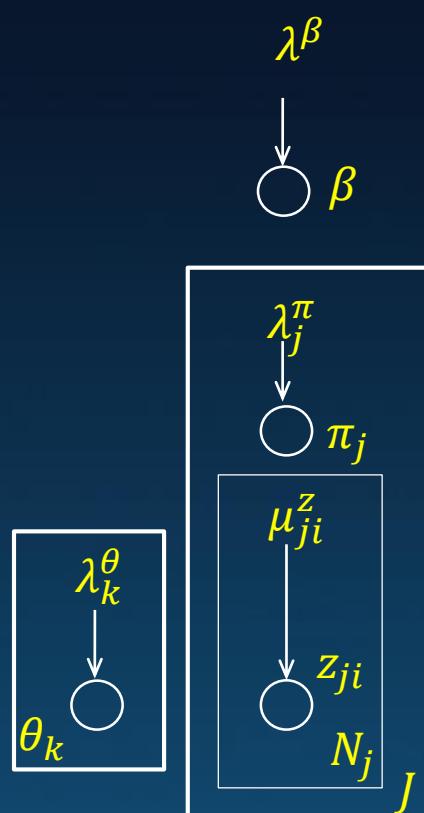
Stochastic Variational Inference (SVI)

Suy diễn biến phân ngẫu nhiên cho mô hình chủ đề

[Hofmann, et al., JMLR 13]

- Giải thuật suy diễn biến phân ngẫu nhiên cho LDA
Loop until convergence

- For $j = 1:J$
 - For $i = 1:N_j$
 - Update $q(z_{ji}) \propto \exp\left(\mathbb{E}_{q-z_{ji}}[\ln p(\Phi, x_{ji})]\right)$
 - Update $q(\pi_j) \propto \exp\left(\mathbb{E}_{q-\pi_j}[\ln p(\Phi, x_j)]\right)$
 - Update $\lambda_k^{\theta(t)} = \lambda_k^{\theta(t-1)} + \rho^{(t)} \frac{|D_t|}{|D|} \frac{\partial F(q(\theta), D_t)}{\partial \lambda_k^\theta}$
 - Update $\lambda^\beta(t) = \lambda^\beta(t-1) + \rho^{(t)} \frac{|D_t|}{|D|} \frac{\partial F(q(\theta), D_t)}{\partial \lambda^\beta(t)}$



Where D_t is a (much smaller and randomly sampled) subset of entire dataset D , while $S = |D_t|$ and $J = |D|$ are sizes of two data sets.

Stochastic Variational Inference (SVI)

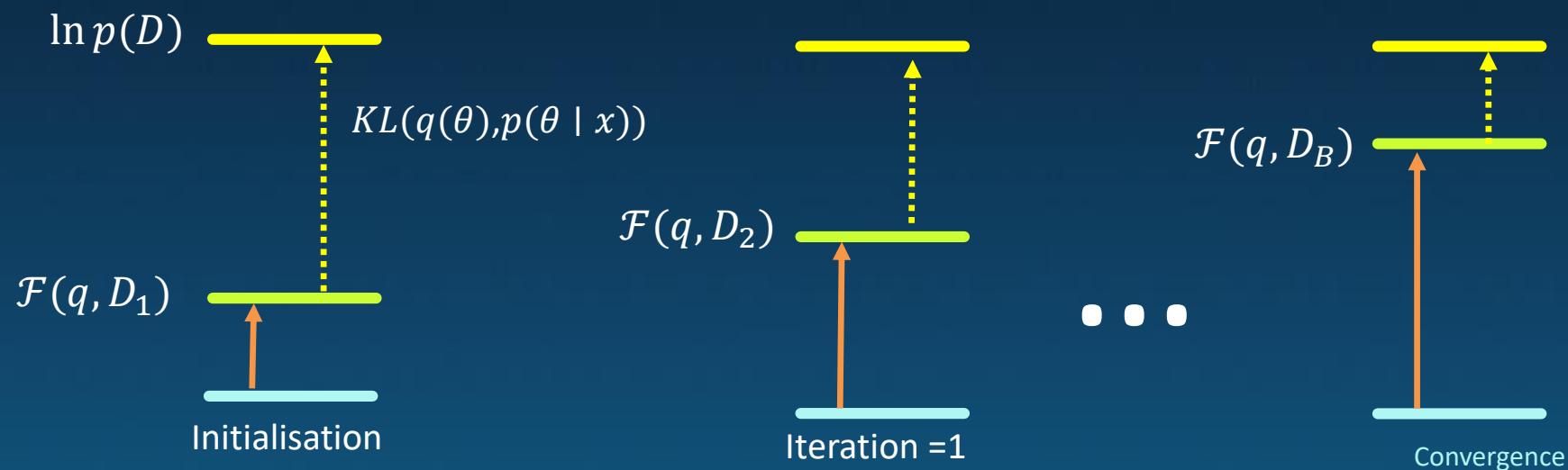
Suy diễn biến phân ngẫu nhiên

- **Variational Inference**

- Each step to update lower bound requires entire data D
 - D is large-scale → **bottleneck**

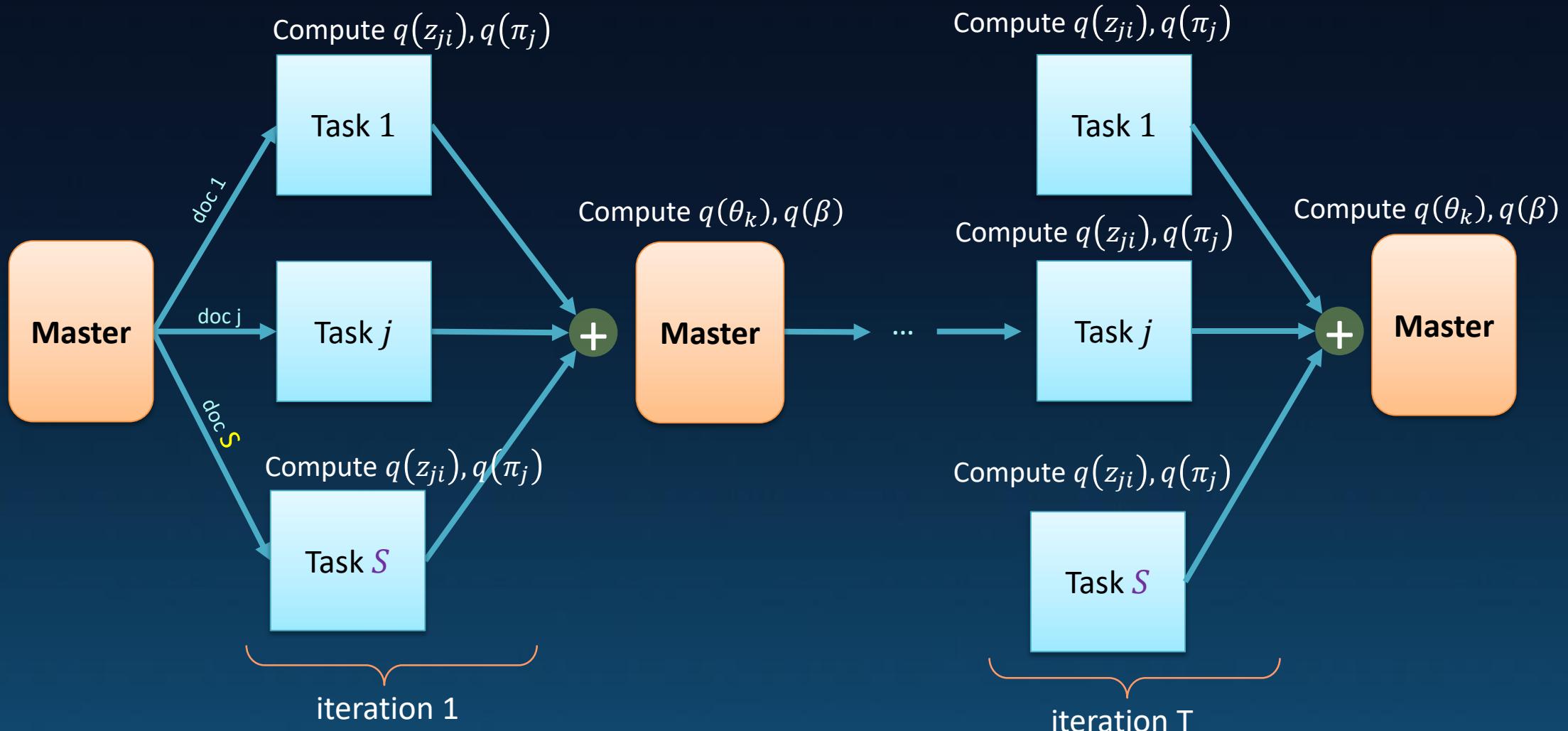
- **Stochastic variational inference**

- use a (randomly chosen small) subset of D , $D_B \subset D$ in each iteration
 - Update the **local parameters** using **coordinate ascent**
 - update the **global parameters** using **stochastic gradient ascent**



Distributed SVI

Xử lý phân tán cho suy diễn biến phân ngẫu nhiên



Now, each iteration involves the computation on much smaller sub dataset D_t ($S \ll J$).

Stochastic Variational Inference (SVI)

Suy diễn biến phân ngẫu nhiên

- Stochastic variational inference

- Lý thuyết graphical models và XSTK cho phép ta xây dựng những mô hình tinh vi (sophisticated models).
- Nhưng những mô hình này (vd: mô hình chủ đề, topic models) không thể áp dụng vào dữ liệu lớn vì độ phức tạp suy diễn quá lớn.
- SVI là một phương pháp cực kỳ hiệu quả giúp chúng ta vượt qua trở ngại này.
- SVI có thể áp dụng vào hầu hết các loại graphical models.

- Ví dụ: làm thế nào để mở rộng mô hình chủ đề (LDA) để không chỉ tìm ra chủ đề ẩn, mà đồng thời phân cụm các tập dữ liệu?

- Huynh, Phung, Nguyen, Hofman and Bui, *Scalable Nonparametric Bayesian Multilevel Clustering*, UAI'16.

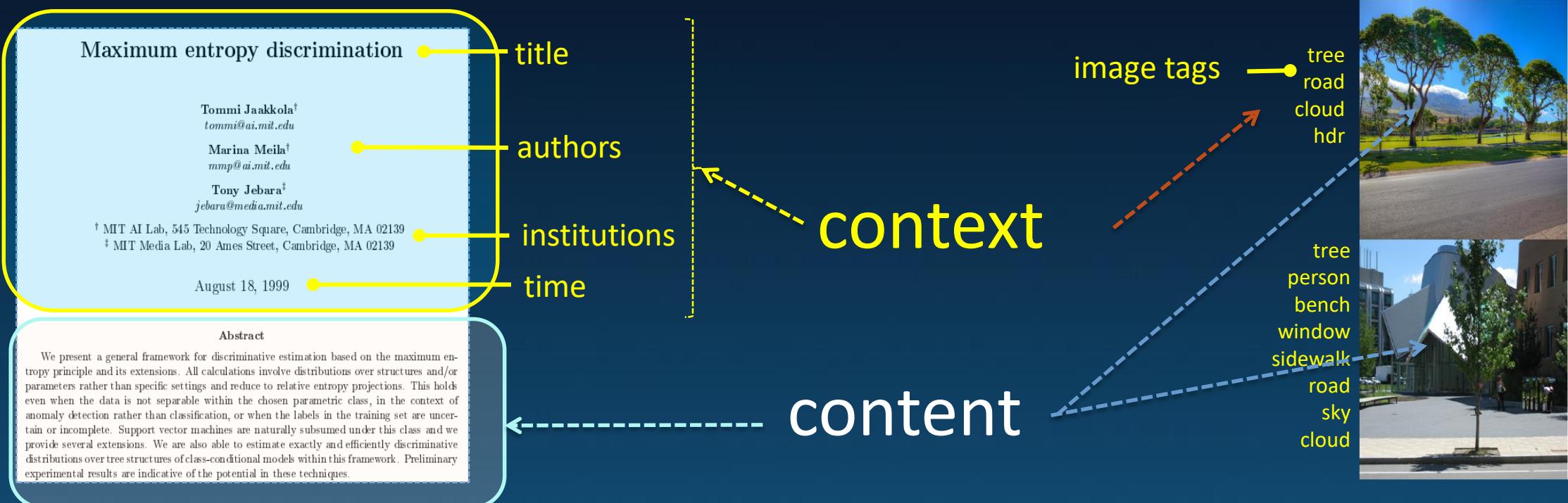
Scalable Bayesian Multilevel Clustering

Gom cụm đa cấp với mô hình chủ đề cho DLL

[Huynh, et al, 16]

• Multilevel clustering with context (MC2)

- Clustering content-units (words, pixels) into “topics”
- Forming content-groups (papers, images) into clusters of data groups
 - Leveraging context-units (title, authors, tags, etc.) associated with each data groups



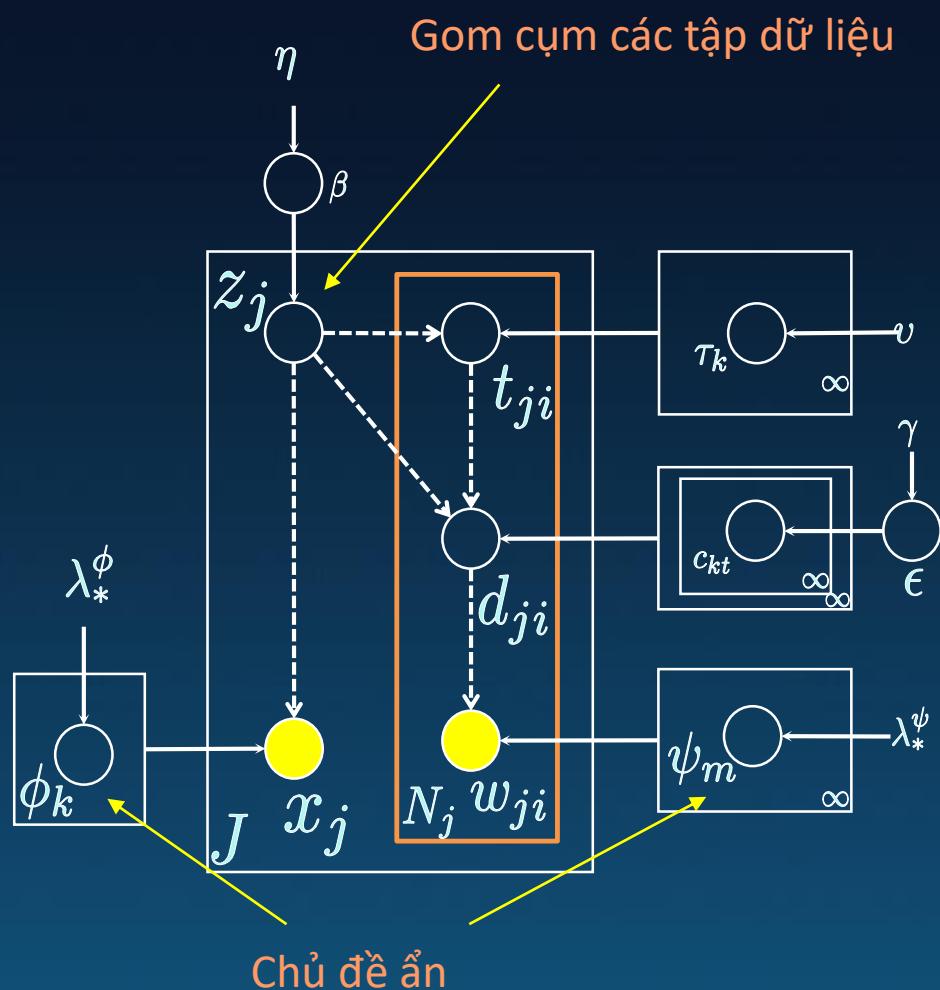
Scalable Bayesian Multilevel Clustering

Gom cụm đa cấp với mô hình chủ đề cho DLL

[Huynh, et al, 16]

- Data $D = \{x_j, w_{ji}\}$
- Parameters $\Theta = \{\beta, z, t, d, c, \tau, \psi, \phi\}$
- Inference:
 - Computing the posterior $p(\Theta | D)$
 - $p(\Theta | D) = \frac{p(D|\Theta)}{\int p(D,\Theta)d\Theta}$ (intractable)
- Phương pháp suy diễn
 - Stochastic Variational Inference (SVI)

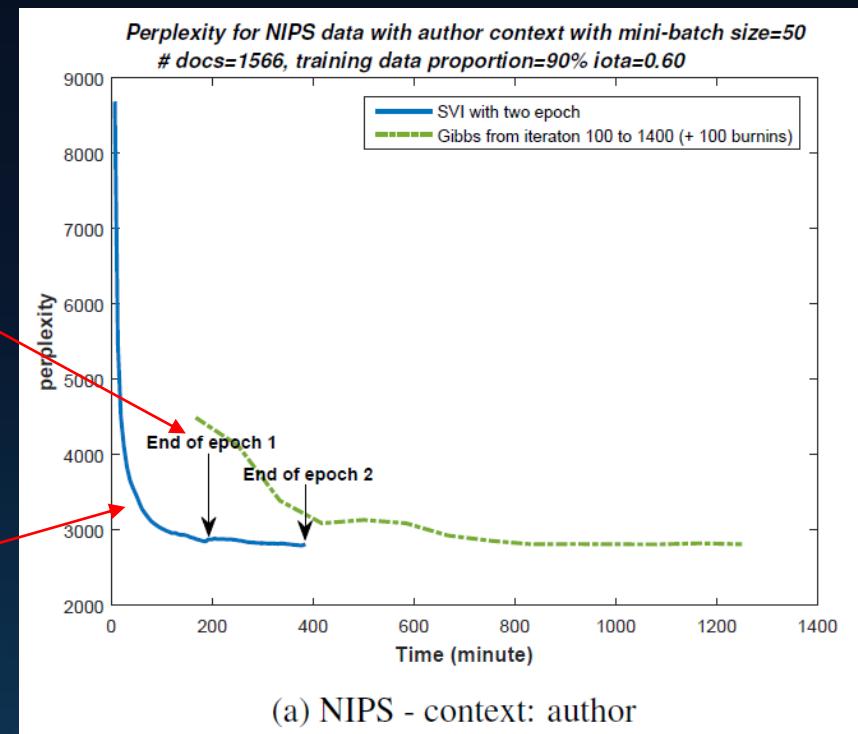
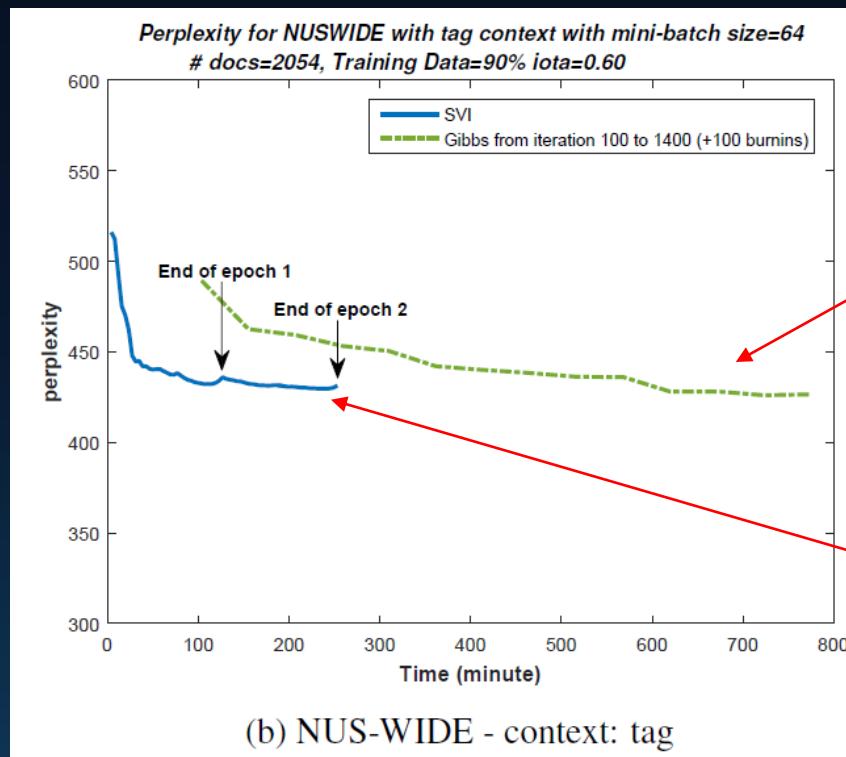
Mô hình MC2



Scalable Bayesian Multilevel Clustering

Tăng tốc giải thuật suy diễn cho MC2 dùng SVI và xử lý song song

[Huynh, et al, 16]



NIPS dataset

NIPS Conference Papers Volumes 0-12 (1987-1999)

1740 papers, 13649 vocabularies

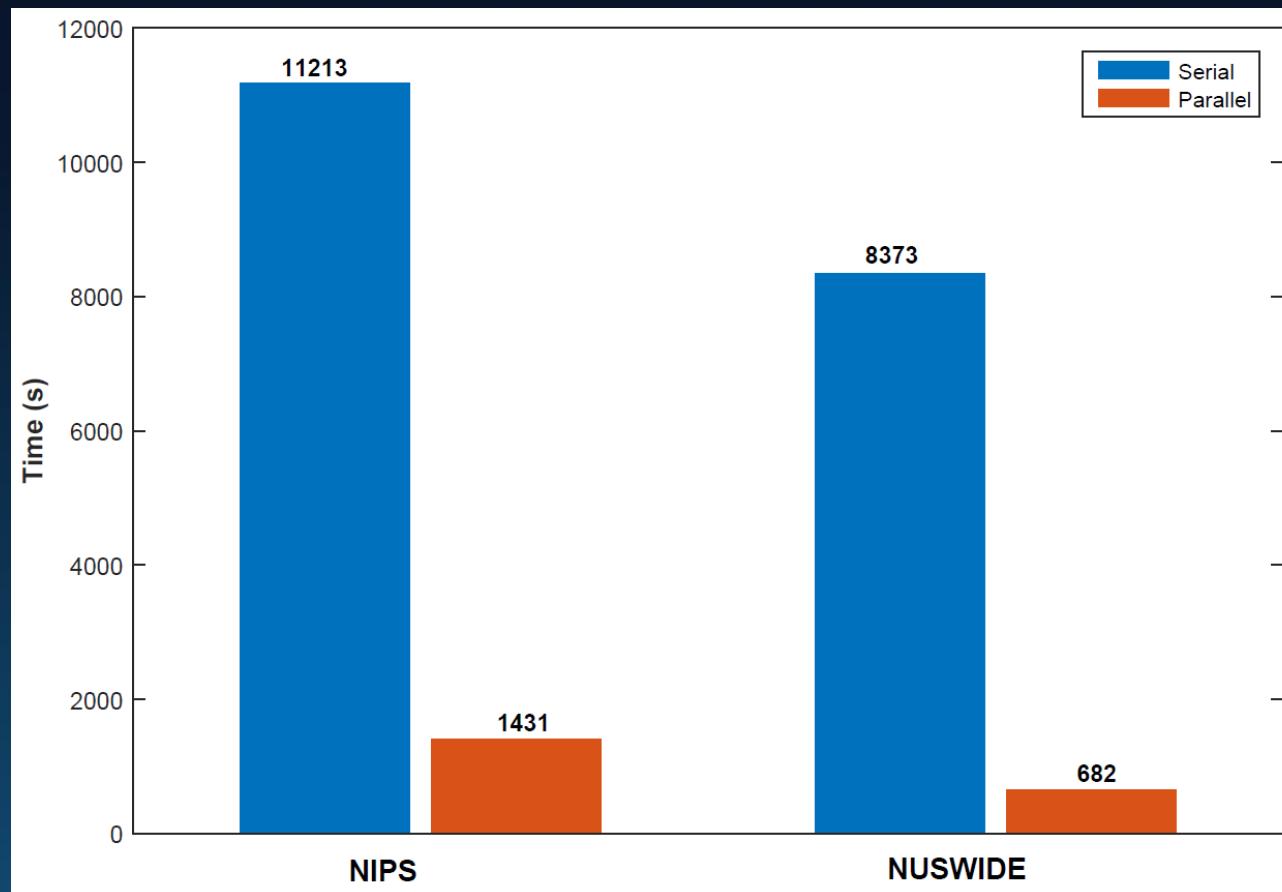
Context information: Authors, Institutions, Published year, Title

Scalable Bayesian Multilevel Clustering

Tăng tốc giải thuật suy diễn cho MC2 dùng SVI và xử lý song song

[Huynh, et al, 16]

Parallel vs sequential implementation of SVI



Scalable Bayesian Multilevel Clustering

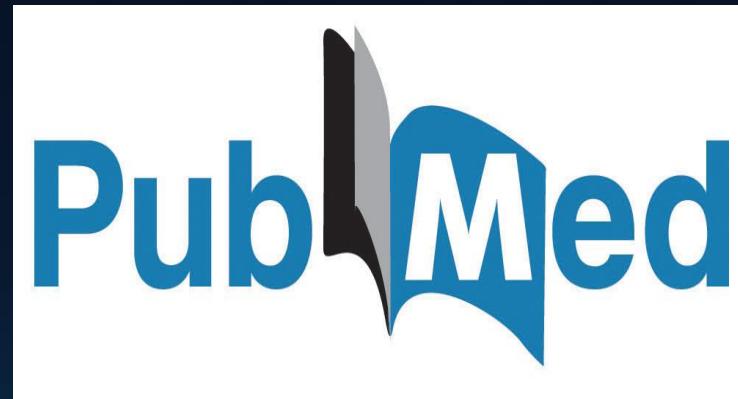
Áp dụng cho các tập dữ liệu lớn

[Huynh, et al, 16]



WIKIPEDIA
The Free Encyclopedia

- 1.1 million documents (billions of data points)
- Vocabulary size 10K
- Contexts:
 - first writer
 - top categories



- 1.4 million abstracts (billions of data points)
- Vocabulary size 10K
- Context:
 - Medical Subject Headings



Application Usage Activity

- >1.1 million users (billions of data points)
- ~ 10K functionalities
- Context:
 - applications
 - paid/trial users

Inference: each iteration update billions of (distributions) of variables !!!

Scalable Bayesian Multilevel Clustering

Kết quả với tập dữ liệu PubMed

[Huynh, et al, 16]

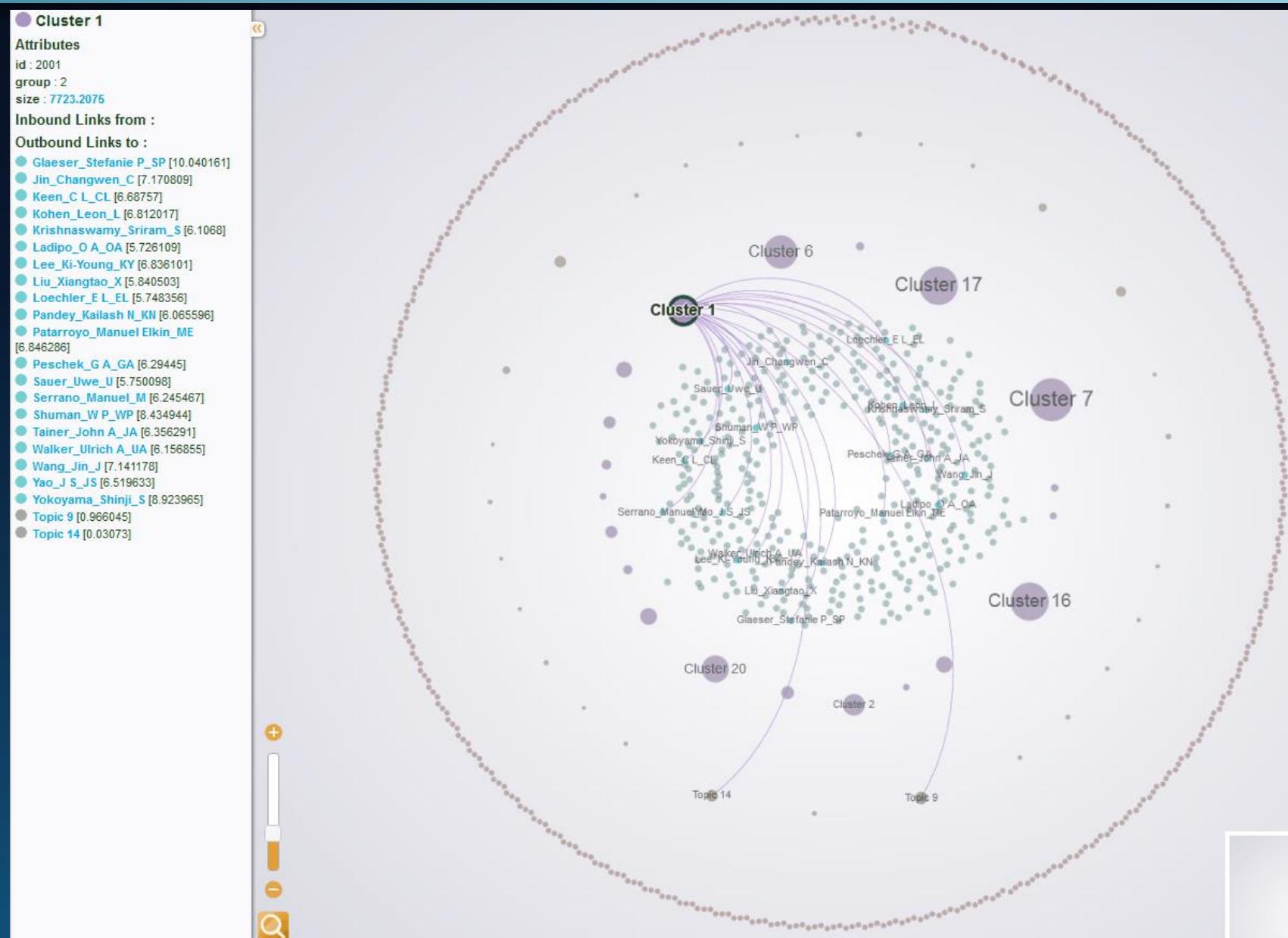
Nodes

- Outermost ring: top **vocabularies** in the corpus.
- Second outermost ring: **topics** discovered by MC2 model.
- Third outermost ring: **cluster of documents** discovered by MC2 model.
- Innermost ring: top **authors** of the papers in the corpus.

Edges between

- vocabularies and topics:** corresponding vocabularies in topics where edge weights denote words' contribution to topics
- clusters and topics:** topics contributes in each cluster of documents
- clusters and authors:** authors contribute clusters of documents

When clicking on each node, e.g. **cluster 1**, the left panel shows nodes connected to active node and normalized weight



Stochastic Variational Inference

Tài liệu tham khảo thêm và mã nguồn

- More on (stochastic) variational inference for graphical models from our work:
 - Huynh, V., Phung, D., Svetha, V., Nguyen, X.L, Hoffman, M. and Bui, H., *Scalable Nonparametric Bayesian Multilevel Clustering*, UAI'16.
 - Huynh, Viet and Phung, Dinh, *Streaming Clustering with Bayesian Nonparametric Models*, Neurocomputing, 2017.
 - Nguyen, V., Phung, D., Venkatesh, S. Nguyen, X.L. and Bui, H, *Bayesian Nonparametric Multilevel Clustering with Group-Level Contexts (MC2)*, ICML'14.
- Source code:
 - <https://github.com/VietHHuynh/MC2SVI>
- PubMed visualization
 - <https://cdn.rawgit.com/VietHHuynh/MC2SVI/7fb53550/visualization/index.html>

Online learning

Ba phương pháp
phát triển mô hình lớn



- **Online learning là gì?**

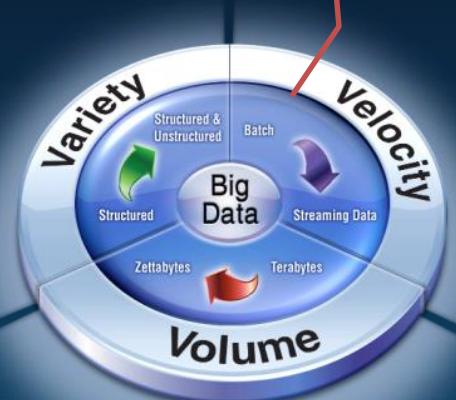
- Trong các hệ thống máy hiện đại, dữ liệu đến hệ thống một cách liên tục.
- Làm thế nào để phát triển mô hình học cập nhật một cách liên tục?
- Online learning: họ các mô hình học gia tăng dần từ dãy dữ liệu nhận được.

- Tại thời điểm $(t - 1)$, thu được tập dữ liệu D_{t-1} và gọi θ_t là tham số mô hình suy diễn

- Tại thời điểm t , tiếp nhận dữ liệu mới (x_t, y_t)
- Cập nhật mô hình theo dữ liệu mới:

$$\theta_{t+1} = h(\theta_t, x_t, y_t)$$

Dòng dữ liệu đến
liên tục và không
ngừng chuyển động



Online learning

- Online learning

- Làm thế nào để đánh giá chất lượng của mô hình học theo cách online?
- Answer: **Regret** represents the retrospect suffered from the online learning algorithm if the optimal model is used to predict the processed data

$$\text{Regret}(\mathbf{u}, T) = \sum_{t=1}^T (l_t(\mathbf{w}_t) - l_t(\mathbf{u}))$$

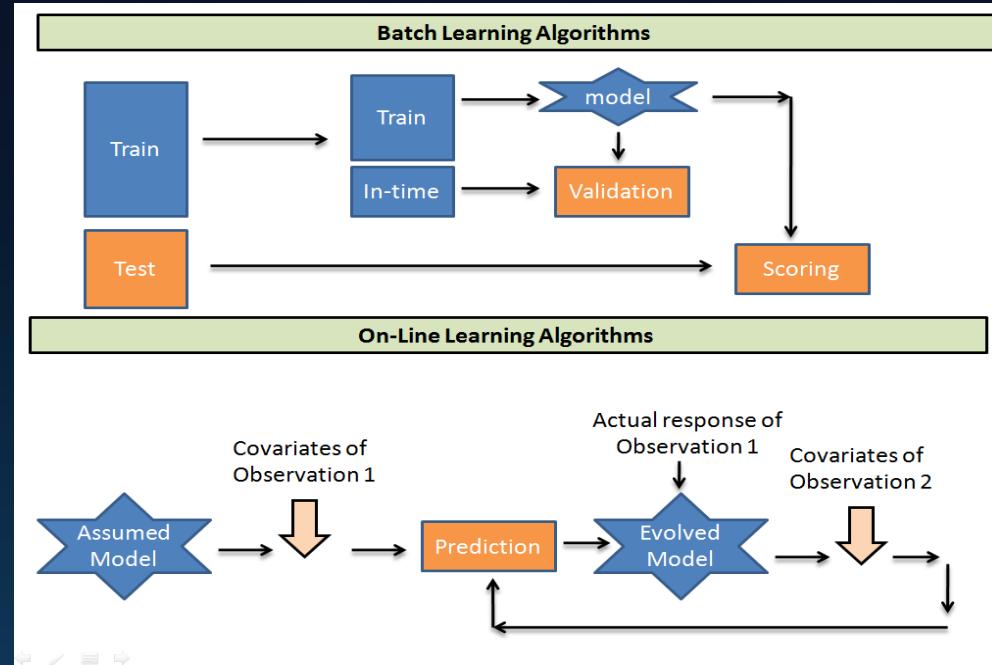
T : số lần cập nhật online
 l_t : loss function (e.g., prediction error for point (x_t, y_t))

$$\text{Regret}(T) = \sup_{\mathbf{u}} \left[\sum_{t=1}^T (l_t(\mathbf{w}_t) - l_t(\mathbf{u})) \right]$$

Online learning

• Online learning

- Giá trị trung bình của regret có thể bị chặn về mặt lý thuyết.
- **Những tốc độ hội tụ lý tưởng:**
 - $\frac{\text{Regret}(T)}{T} \in O\left(\frac{1}{\sqrt{T}}\right), O\left(\frac{\log T}{T}\right), O\left(\frac{1}{T}\right)$
 - Như vậy, nếu khéo léo ta có thể thiết kế mô hình bảo đảm tính hội tụ đến lời giải tốt ưu với tốc độ hội tụ nhanh khi T tăng dần.



Kernel Online learning

Ba phương pháp phát triển mô hình lớn

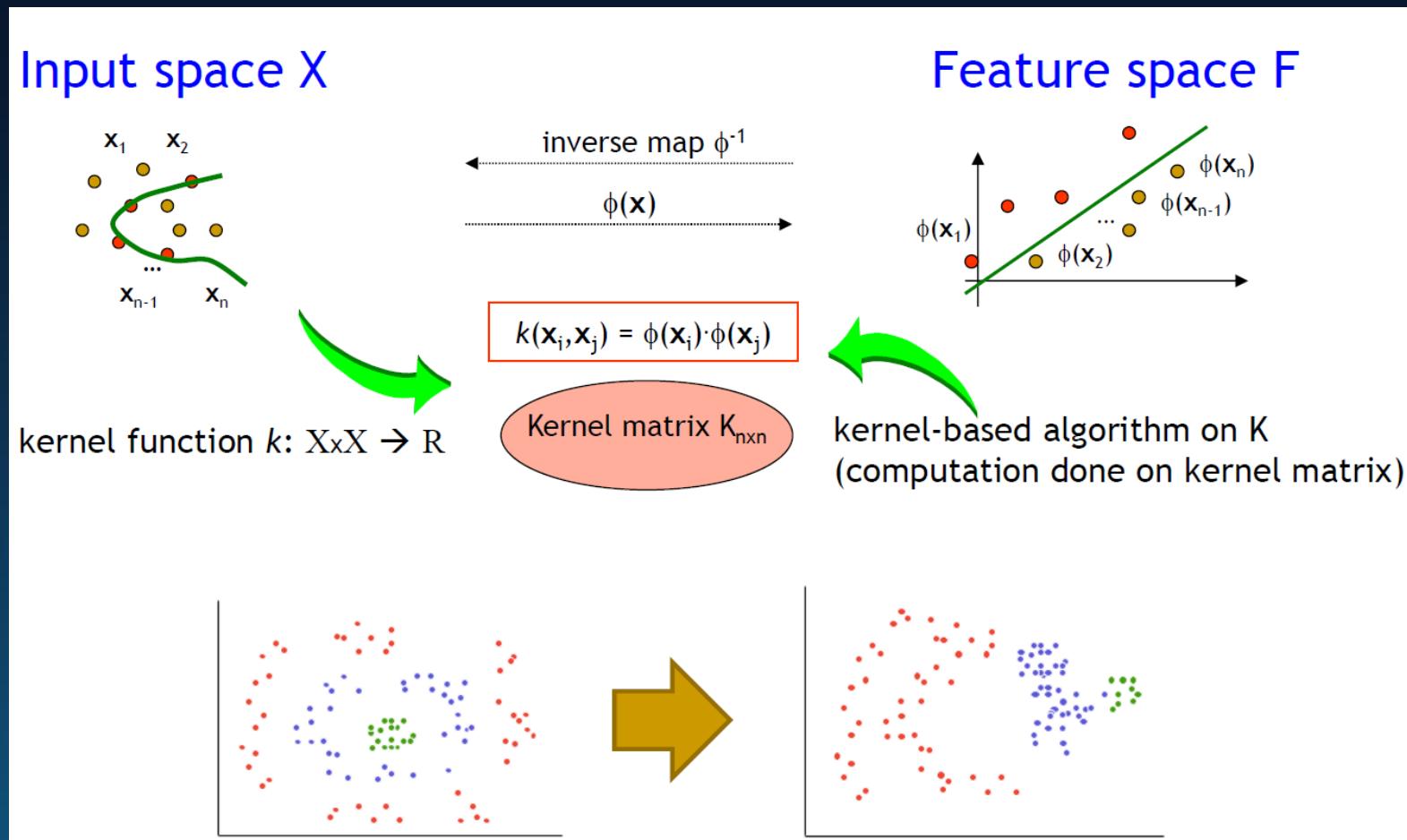
Data Augmentation

Stochastic Variational Inference

Online Learning

• Kết hợp phương pháp kernel và online learning

Kernel trick: ta chỉ cần định nghĩa kernel function mà không cần chỉ ra dạng tương minh của hàm chiếu $\phi(\cdot)$



Kernel online learning

- Kết hợp thế mạnh của mô hình kernel và online learning.
- Typically, solution for a kernel method applied to online case:
 - Current model $\theta_t = \sum_{i=1}^t \alpha_i \phi(x_i)$ where $\phi(\cdot)$ is feature map
 - Model size is defined as $\|\alpha\|_0$
 - Problem: curse of kernelization, i.e., model size increases linearly over time.
 - $\phi(\cdot)$ lies in infinite-space, hence we can't store it directly – all computation is done via the kernel function and α .
- Two approaches to combat the curse of kernelization
 - Budgeted [Wang et al 2012]
 - Limit model size to a predefined budget B (removal, merging, projection)
 - Random features [Lu et al 2016]
 - Express $\phi(\cdot)$ explicitly via Bochner's theorem

Kernel online learning

Bochner's theorem [Bochner, 1959]: *every shift-invariant kernel $k(x, x')$ can be represented as an inverse Fourier transform of a proper distribution $p(\omega)$ as below:*

$$k(x, x') = k(u) = \int p(\omega) e^{i\omega^\top u} du \text{ where } u = x - x'$$

- Consequently, given the form of the kernel, one can derive an explicit distribution for ω : $p(\omega) = \left(\frac{1}{2\pi}\right)^d \int k(u) e^{-iu^\top \omega} du$
 - E.g, if the kernel is a Gaussian kernel, then ω follows a Gaussian distribution.

Using Euler's formula, one can approximate: $k(x, x') \approx \tilde{k}(x, x') = \tilde{\phi}^\top(x)\tilde{\phi}(x')$

where

$$\tilde{\phi}(x) = \sqrt{\frac{2}{D}} [\cos(\omega_j^\top x), \sin(\omega_j^\top x)]_{j=1}^D$$

and

$$\omega_j \stackrel{iid}{\sim} p(\omega), j = 1, 2, \dots, D$$

Kernel online learning

Bochner's theorem [Bochner, 1959]: every shift-invariant kernel $k(x, x')$ can be represented as an inverse Fourier transform of a proper distribution $p(\omega)$ as below:

$$k(x, x') = k(u) = \int p(\omega) e^{i\omega^\top u} du \text{ where } u = x - x'$$

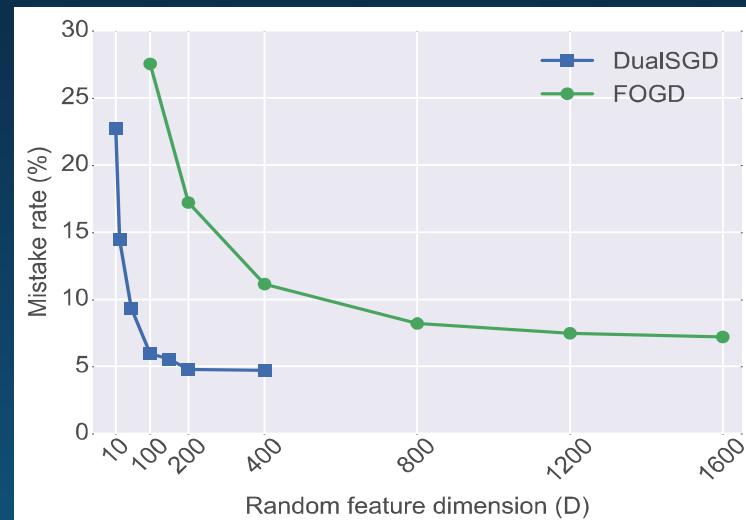
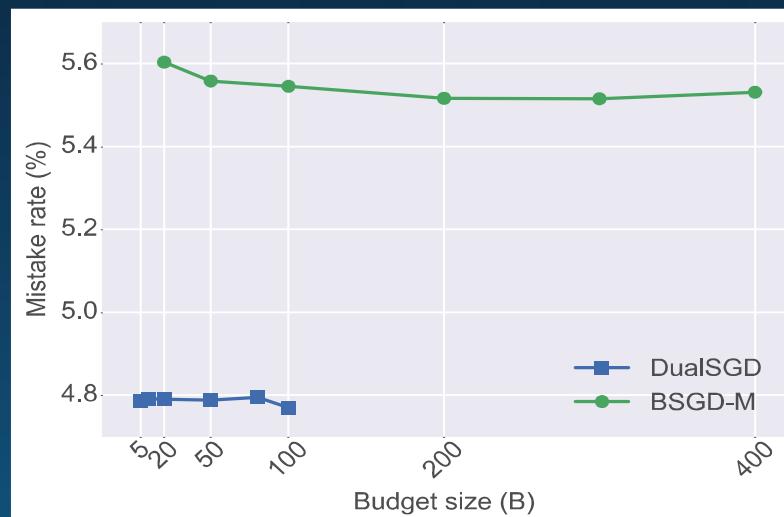
- Định lý Bochner cho phép ta xây dựng hàm chiểu $\phi(\cdot)$ một cách tương minh.
- Kết quả này giúp chúng ta xây dựng nhiều mô hình mới!
 - Large-scale Online Kernel Learning with Random Feature Reparameterization* [Tu et al., IJCAI 17]
 - Kernel Online learning for both model parameters and kernel parameters.
 - Dual Space Gradient Descent for Online Learning* [Le et al, NIPS16]
 - Advanced kernel online learning method that ensures theoretical convergence and quality.

Gradient Descent trên không gian kép

Dual Space Gradient Descent for Online Learning

[Le et al., NIPS16]

- Random feature approach [Lu et al 2015] operates on *random-feature space* (using Bochner's theorem), and then perform SGD in the feature space.
- Problem: excessive number of random features is required
- DualSGD [Le et al. NIPS16]: use both original feature space and random feature space
 - Key intuition: use *a provision vector* in the random-feature space to *store* the information of all vectors being removed via *k-merging* technique.
 - Guarantee in terms of theory for convergence and accuracy.



Tham khảo

References

- For more online learning papers from our group, see:
 - Trung Le, Tu Nguyen, Vu Nguyen and Phung, Dinh, *Approximation Vector Machines*, Journal of Machine Learning Research (**JMLR**, 2017)
 - Trung Le, Tu Nguyen, Vu Nguyen and Dinh Phung, *Dual Space Gradient Descent for Online Learning*, **NIPS'16**.
 - Trung Le, Vu Nguyen, Tu Nguyen and Dinh Phung, *Nonparametric Budgeted Stochastic Gradient Descent*, **AISTATS'16**.
 - Trung Le, Duong Phuong, Dinh Mi, Tu Nguyen and Dinh Phung, *Budgeted Semi-supervised Support Vector Machine*, **UAI'16**.
 - Tu Nguyen, Trung Le, Hung Bui and Dinh Phung, *Large-scale Online Kernel Learning with Random Feature Reparameterization*, **IJCAI'17**.

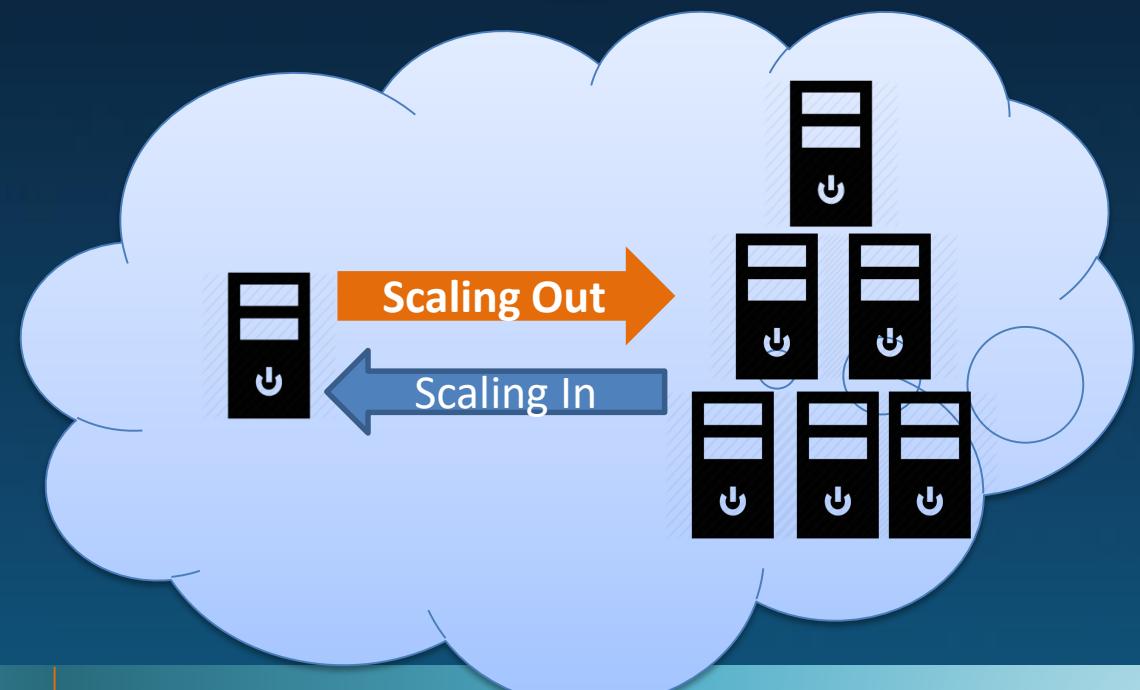
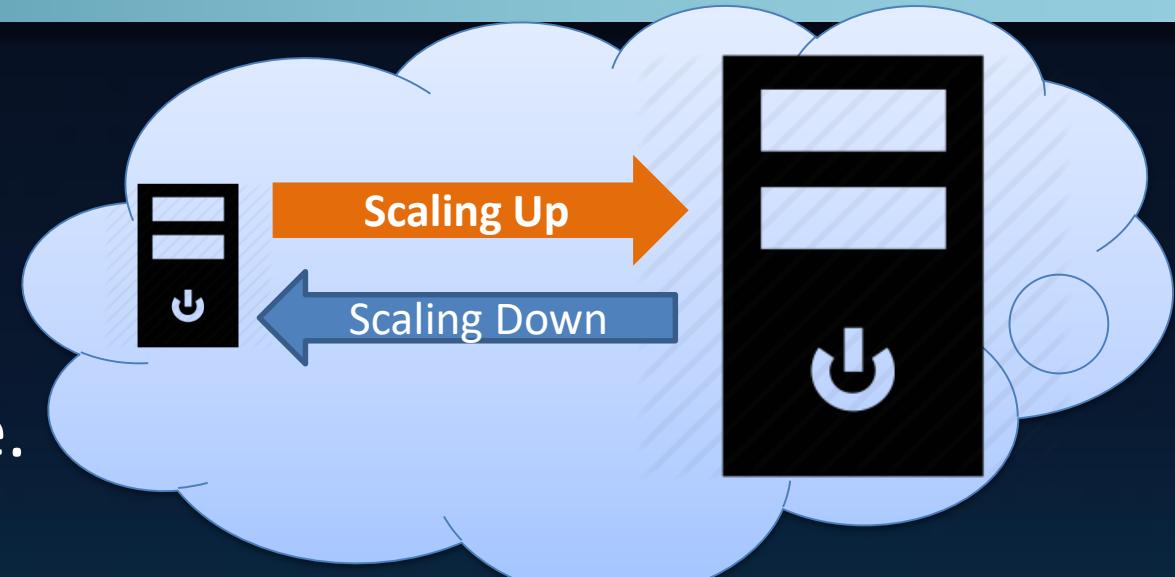
Công nghệ nào cho mô hình lớn?

Platform for big models

- Mô hình lớn cho dữ liệu lớn
 - Three approaches to scale up big models
 - Scale up Gradient Descent Methods
 - Scale up Probabilistic Inference
 - Scale up Model Evaluation (bootstrap)
- Ba cách tiếp cận để phát triển mô hình lớn
 - Data augmentation
 - Stochastic Variational Inference for Graphical Models
 - Stochastic Gradient Descent and Online Learning
- Công nghệ nào cho mô hình lớn?
 - Hadoop, Spark, TensorFlow

Platform for Big Data

- Scaling up (Vertical Scaling): adding more power (CPU, RAM) to an existing machine.
- Scaling out (Horizontal Scaling): adding more machines into the pool of resources.



Platform for Big Data

Ưu và nhược điểm của mở rộng theo chiều ngang và chiều dọc

Horizontal Scaling

Increases performance in **small steps** as needed
Financial investment to upgrade is relatively less
Can scale out the system **as much as needed**



Horizontal Scaling

Software has to handle all the data distribution and parallel processing **complexities**
Limited number of softwares are available that can take advantage of horizontal scaling

Vertical Scaling

Most of the **software** can easily take advantage of vertical scaling
Easy to manage and install hardware within a single machine

Vertical Scaling

Requires **substantial financial investment**
System has to be **more powerful** to handle **future workloads** and initially the additional performance is not fully utilized
It is **not possible** to scale up vertically **after a certain limit**

Platform for Big Data

Horizontal Scaling with Peer-to-Peer Networks

- Peer-to-peer networks

- millions of machines connected in a network
- communication scheme exchange the data:
Message Passing Interface (MPI)
- Advantages:

- well suited for iterative processing
- available for many programming languages
- useful for dynamic resource allocation with the master–slave model

- Disadvantages:

- no mechanism to handle faults
- more complex with lots of functions

Horizontal
Scaling

Công nghệ MapReduce

Parallel Processing with MapReduce



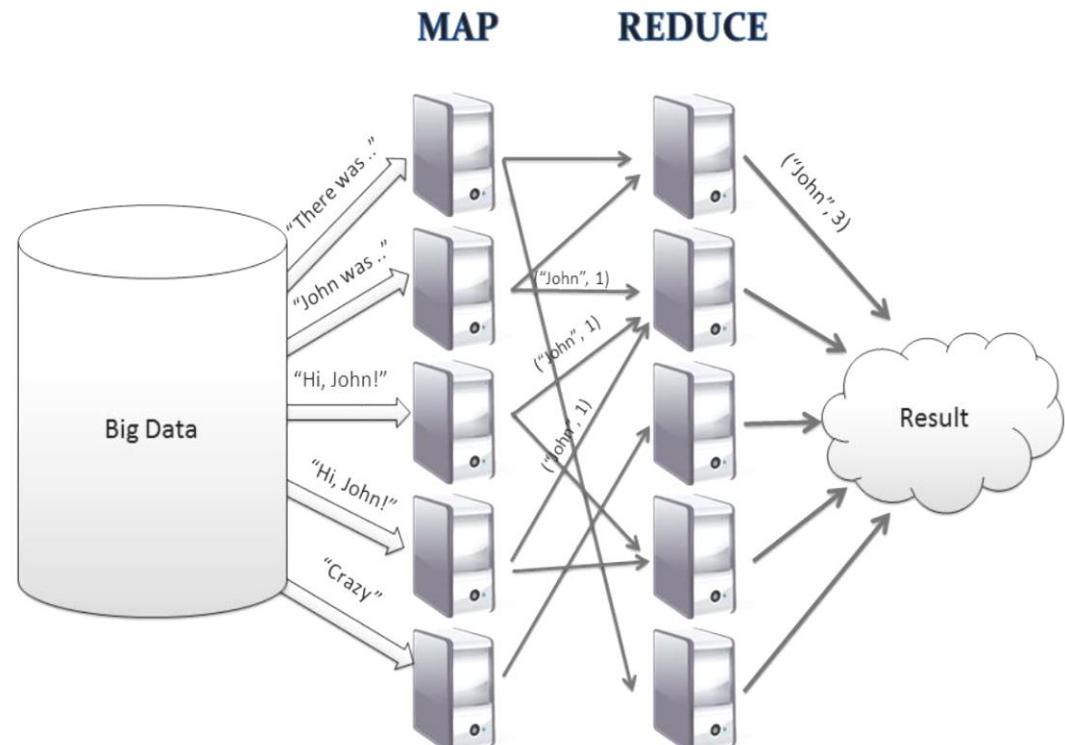
- Nhu cầu xử lý DLL tăng nhanh

- MapReduce

- invented by Google in 2004 and used in Hadoop.
- breaking the entire task into two parts: **mappers** and **reducers**.
- **mappers**: read the data from HDFS, process it and generate some intermediate results.
- **reducers**: aggregate the intermediate results to generate the final output.

- Key Limitations

- inefficiency in running iterative algorithms.
- Mappers read the same data again and again from the disk.



Data parallelization

Distributed execution

Outcome aggregation

Source: <https://www.supinfo.com/articles/single/2807-introduction-to-the-mapreduce-life-cycle>

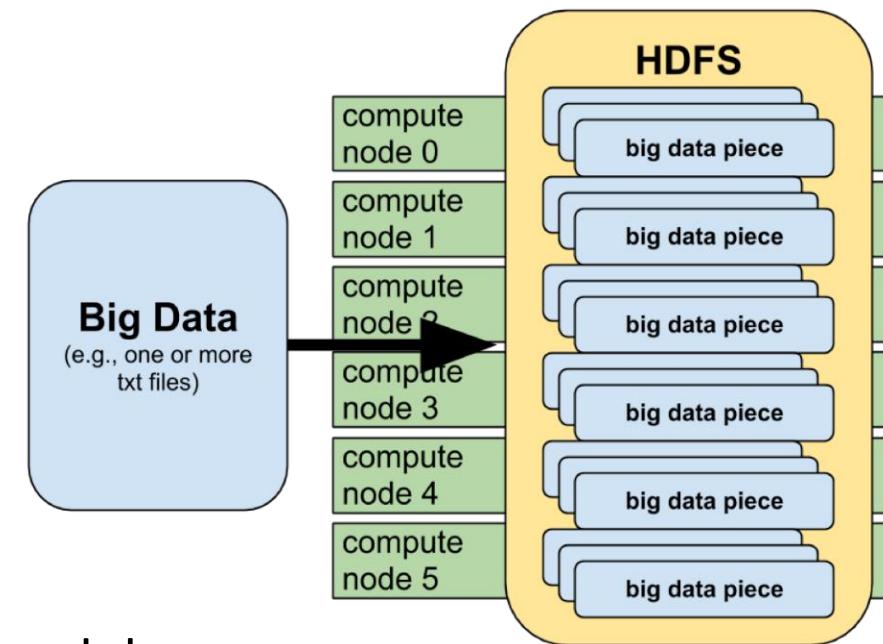
Công nghệ Hadoop

Apache Hadoop



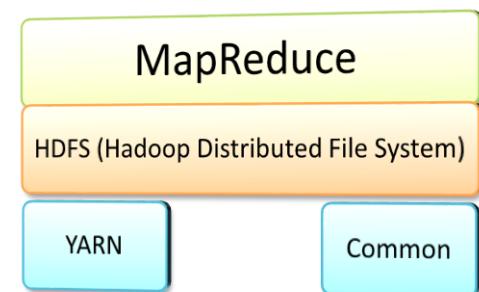
❑ Apache Hadoop:

- an open source framework for storing and processing **large datasets** using clusters of commodity hardware
- **highly fault tolerant**
- **scaling up to 100s or 1000s of nodes**



❑ Hadoop components:

- **Common:** utilities that support the other Hadoop modules
- **YARN:** a framework for job scheduling and cluster resource management.
- **HDFS:** a distributed file system
- **MapReduce:** computation model for parallel processing of large datasets.



❑ **Key limitation:** not suitable for iterative-convergent algorithm!

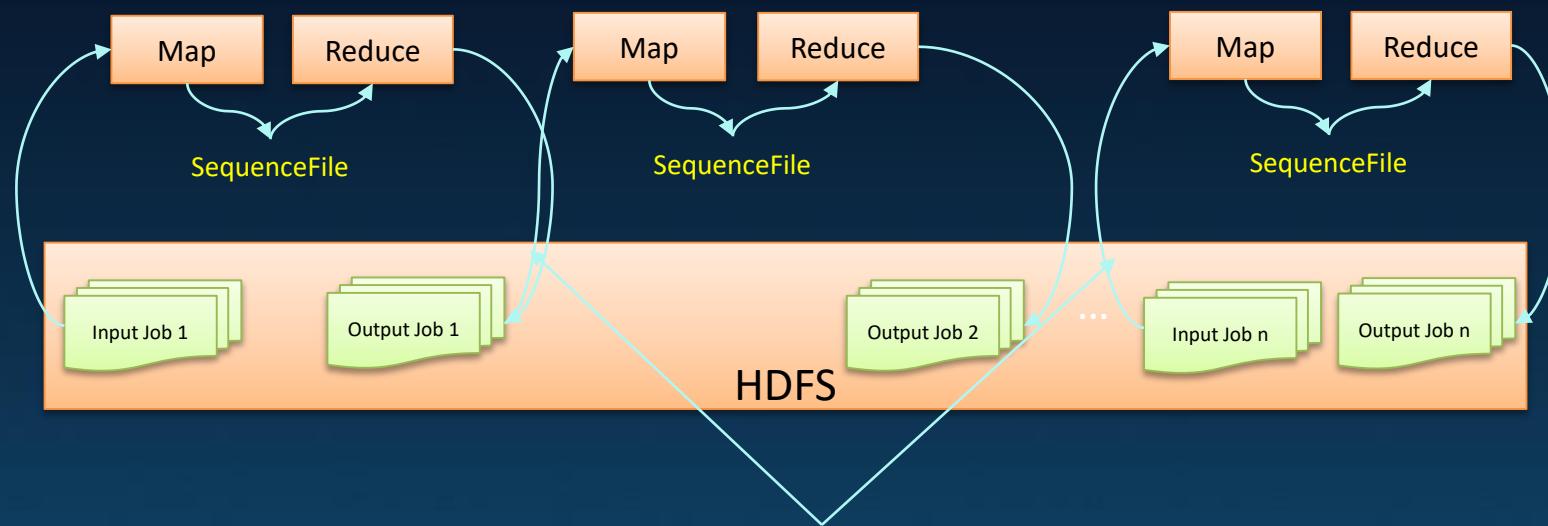
Source: opensource.com

Platform for Big Data

Hadoop có phù hợp với các giải thuật lặp lại?

• Limitations of MapReduce

- ❑ inefficiency in running iterative algorithms
- ❑ Mappers read the same data again and again from the disk



File access overhead!!!

How to reduce this latency? → memory caching

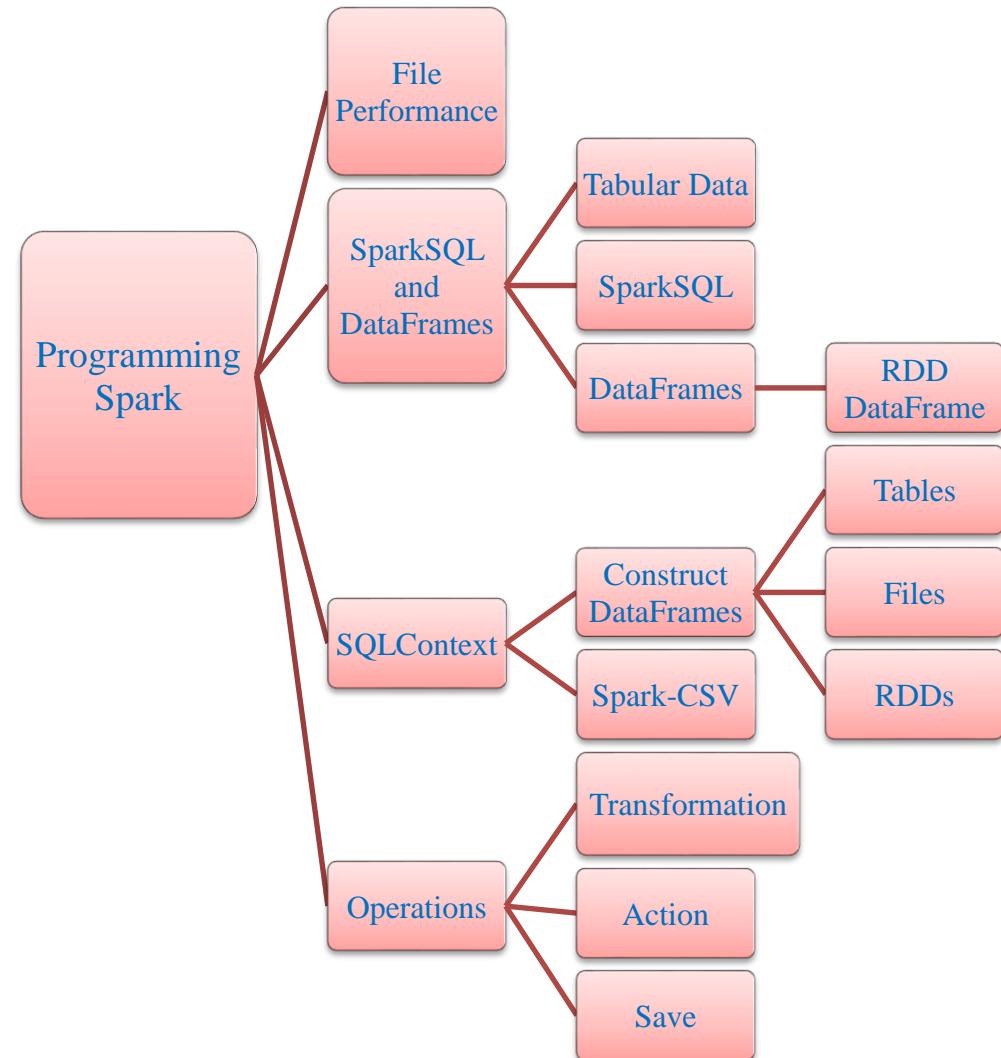
Hadoop có phù hợp với các giải thuật lặp lại? Không thật sự phù hợp

Công nghệ Spark

Apache Spark



- Key motivation: suitable for iterative-convergent algorithms!
- Spark key features:
 - Resilient Distributed Datasets (RDD)
 - Read-only, partitioned collection of records distributed across cluster, stored in memory or disk.
 - Data processing = graph of transforms where nodes = RDDs and edges = transforms.
 - Benefits:
 - Fault tolerant: highly resilient due to RDDs
 - Cacheable: store some RDDs in RAM, hence faster than Hadoop MR for iteration.
 - Support MapReduce.

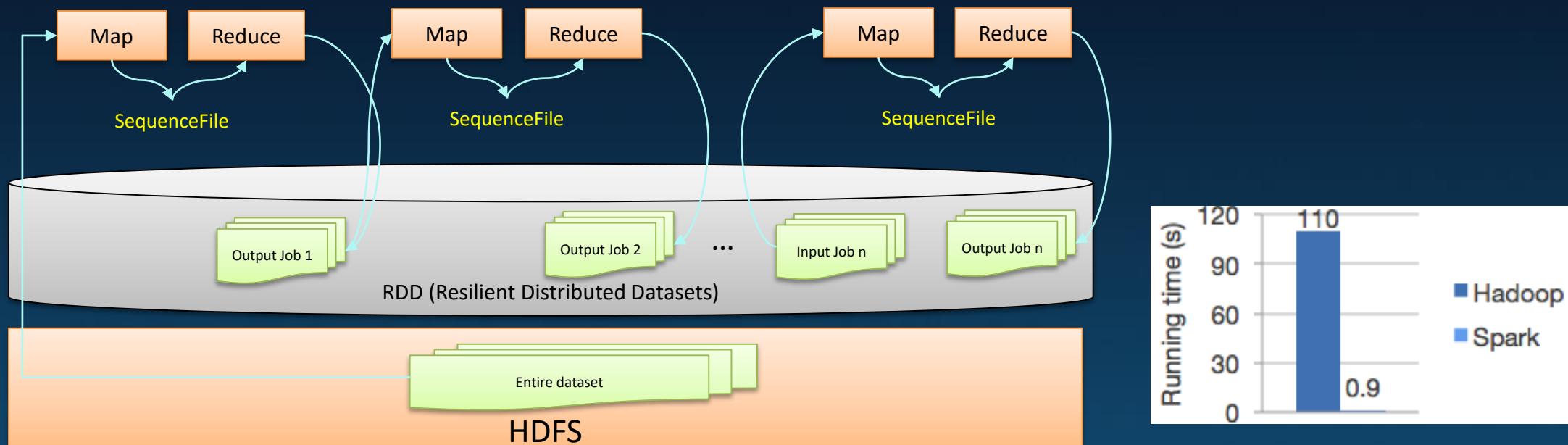


Platform for Big Data

Apache Spark thích hợp cho tính toán lặp lại

Apache Spark

- an alternative to Hadoop and MapReduce
- perform in-memory computations

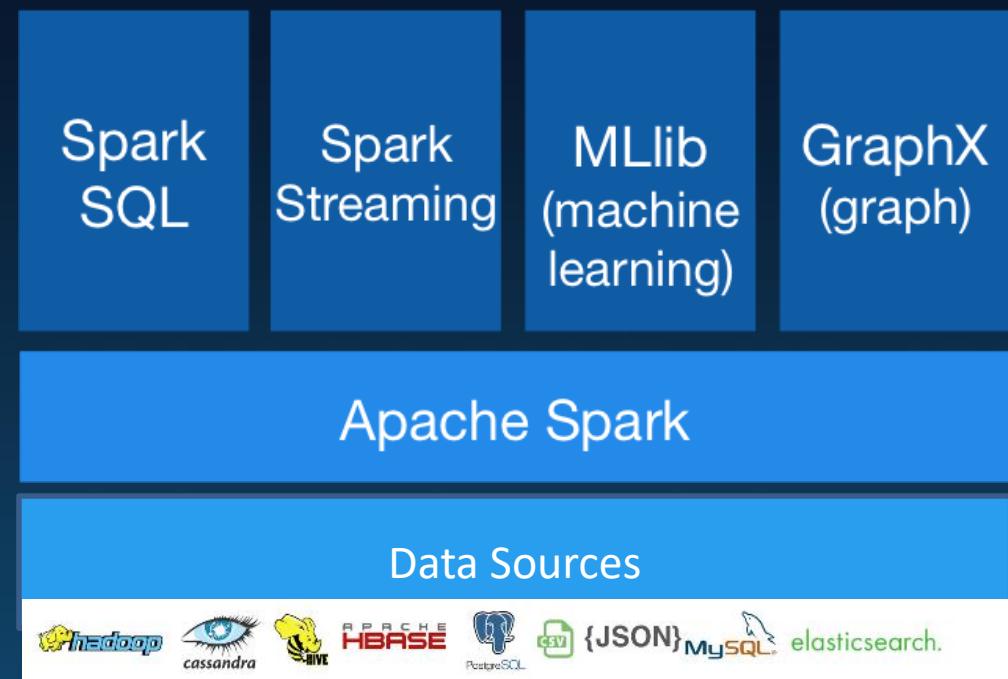


Apache Spark thích hợp cho tính toán lặp lại!

Apache Spark

Các thành phần của Apache Spark

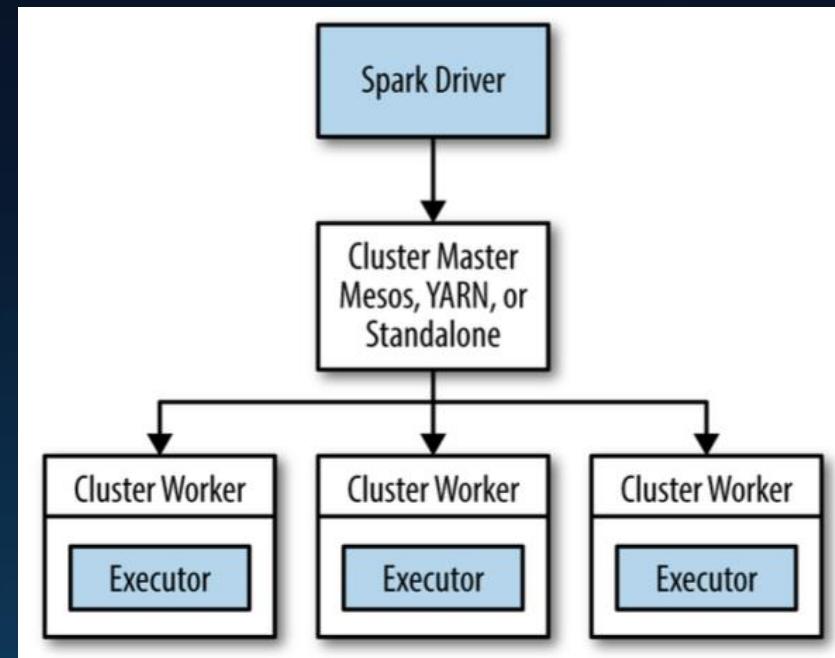
- Spark Core includes
 - components for task scheduling, memory management, fault recovery, interacting with storage systems
 - the API for resilient distributed datasets (RDDs)
- Spark SQL
 - package for working with structured data
- Spark Streaming
 - component for processing of live streams of data
- MLlib
 - library containing common machine learning (ML) functionality
- GraphX
 - library for manipulating graphs and performing graph-parallel computations



Apache Spark

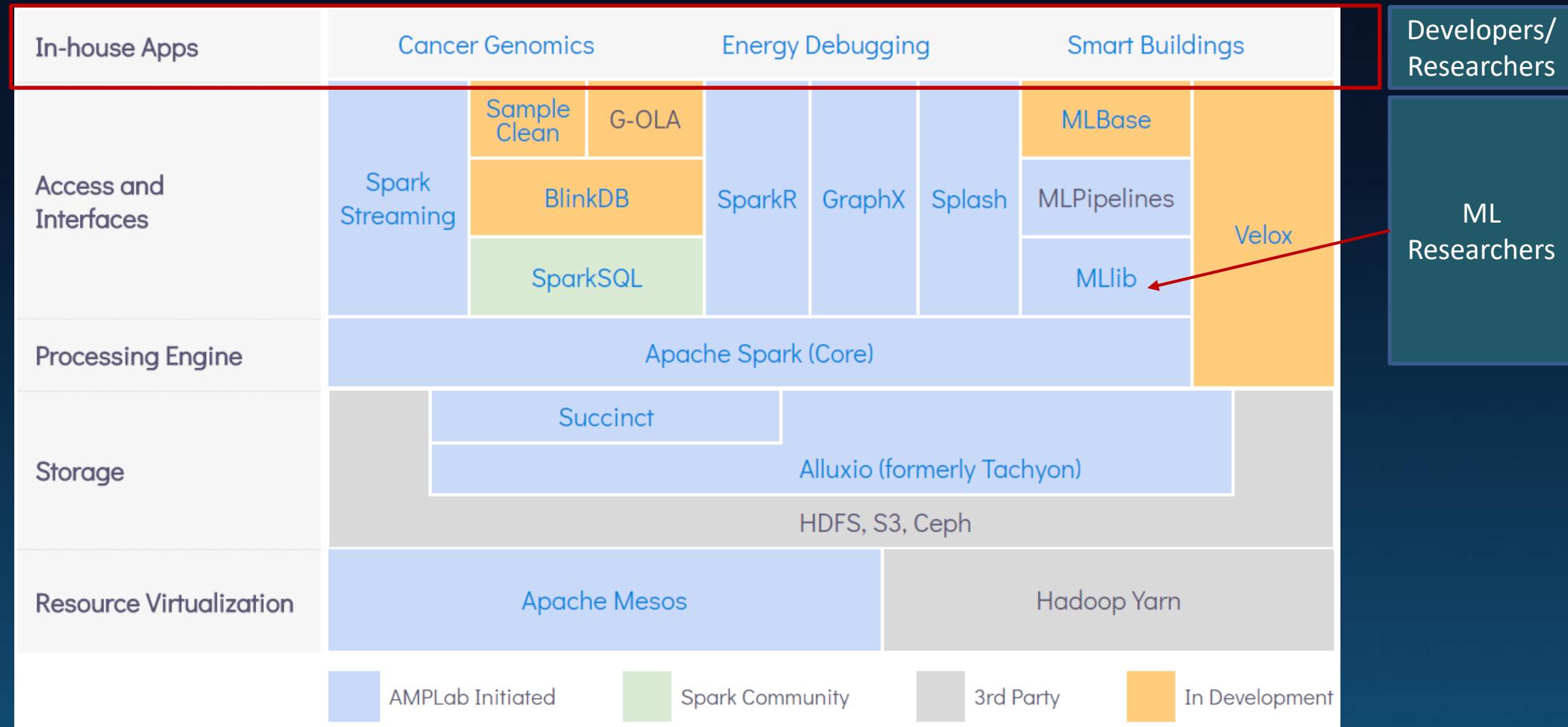
Kiến trúc thực thi của Spark

- In distributed mode, Spark uses a master/slave architecture
 - one central coordinator (**driver**) and
 - many distributed workers (**executors**)
- Spark Application contains
 - The **driver** runs in its own Java process
 - Each **executor** is a separate Java process



Ví dụ: BDAS, Berkeley Data Analytics Stack

Bộ công cụ tích hợp cho phân tích dữ liệu lớn



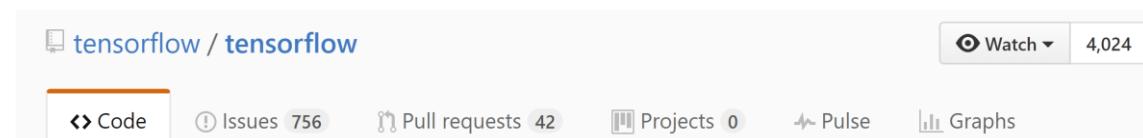
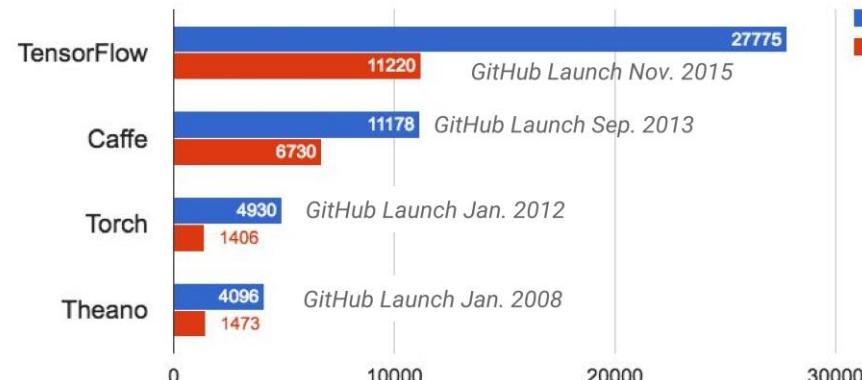
[Nguồn: <https://amplab.cs.berkeley.edu/software/>]

TensorFlow

- open-source framework for deep learning, developed by the GoogleBrain team.
- provides primitives for defining functions on tensors and automatically computing their derivatives.

Strong External Adoption

Adoption of Deep Learning Tools on GitHub



Screenshot of the TensorFlow GitHub repository page. The repository name is "tensorflow / tensorflow". It shows 4,024 stars and 13,142 commits. The repository description is "Computation using data flow graphs for scalable machine learning <http://tensorflow.org>".

Computation using data flow graphs for scalable machine learning <http://tensorflow.org>

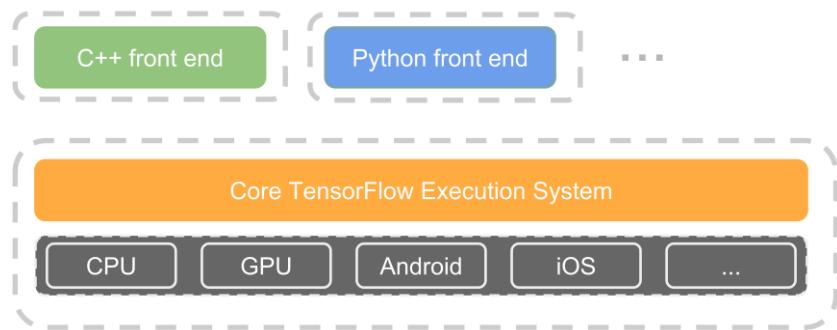
Công nghệ TensorFlow

Parallel Processing with TensorFlow



TensorFlow

- Back-end in C++: very low overhead
- Front-end in Python or C++: friendly programming language



- ... and even in more platforms

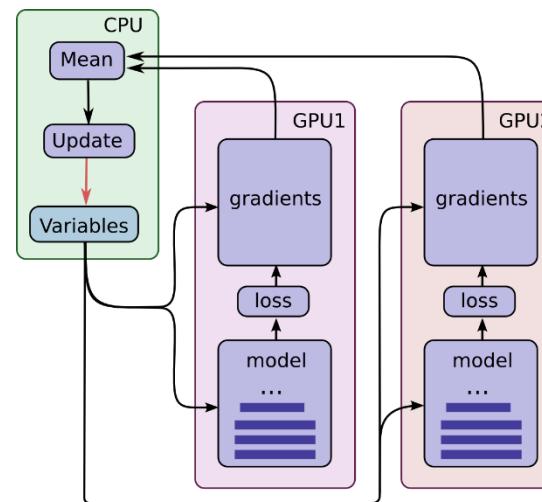
Linux CPU	Linux GPU	Mac OS CPU	Windows CPU	Android
build passing				

- Switchable between CPUs and GPUs

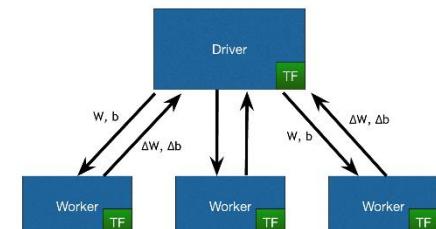
```
import tensorflow as tf  
with tf.device('/cpu:0'):  
    a = tf.constant(1.0)
```

```
import tensorflow as tf  
with tf.device('/gpu:0'):  
    a = tf.constant(1.0)
```

- Multiple GPUs in one machine or distributed over multiple machines



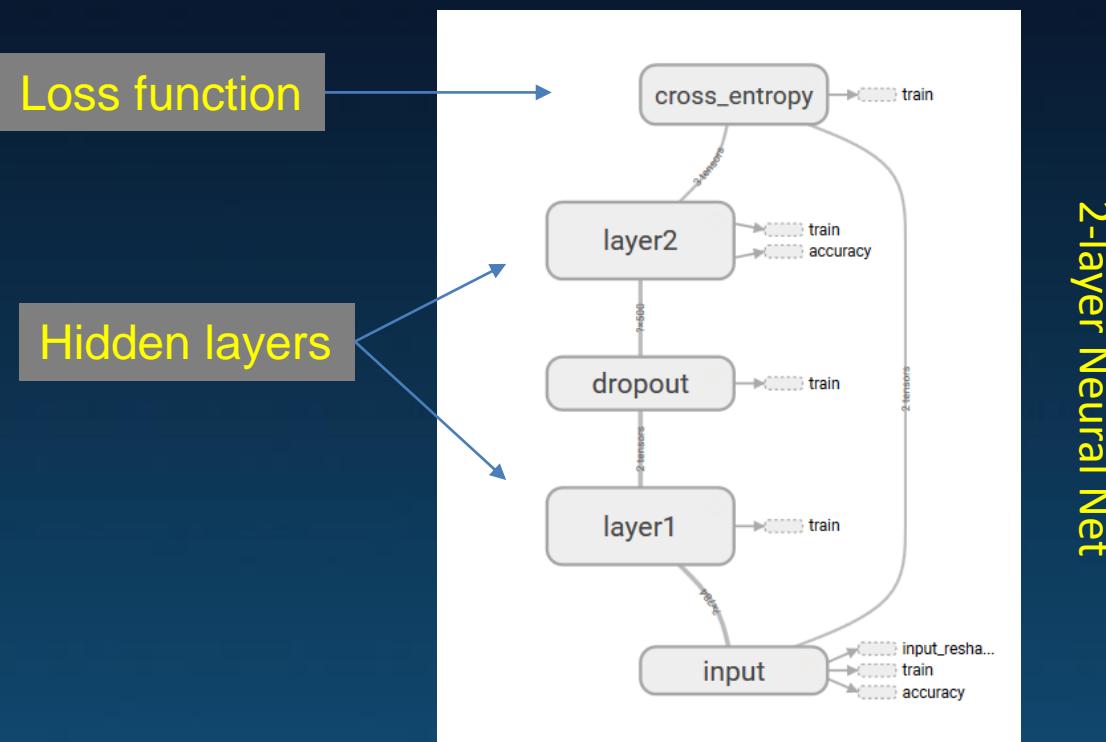
```
import tensorflow as tf  
with tf.device('/gpu:0'):  
    a = tf.constant(1.0)  
with tf.device('/gpu:1'):  
    b = tf.constant(2.1)
```



Công nghệ TensorFlow

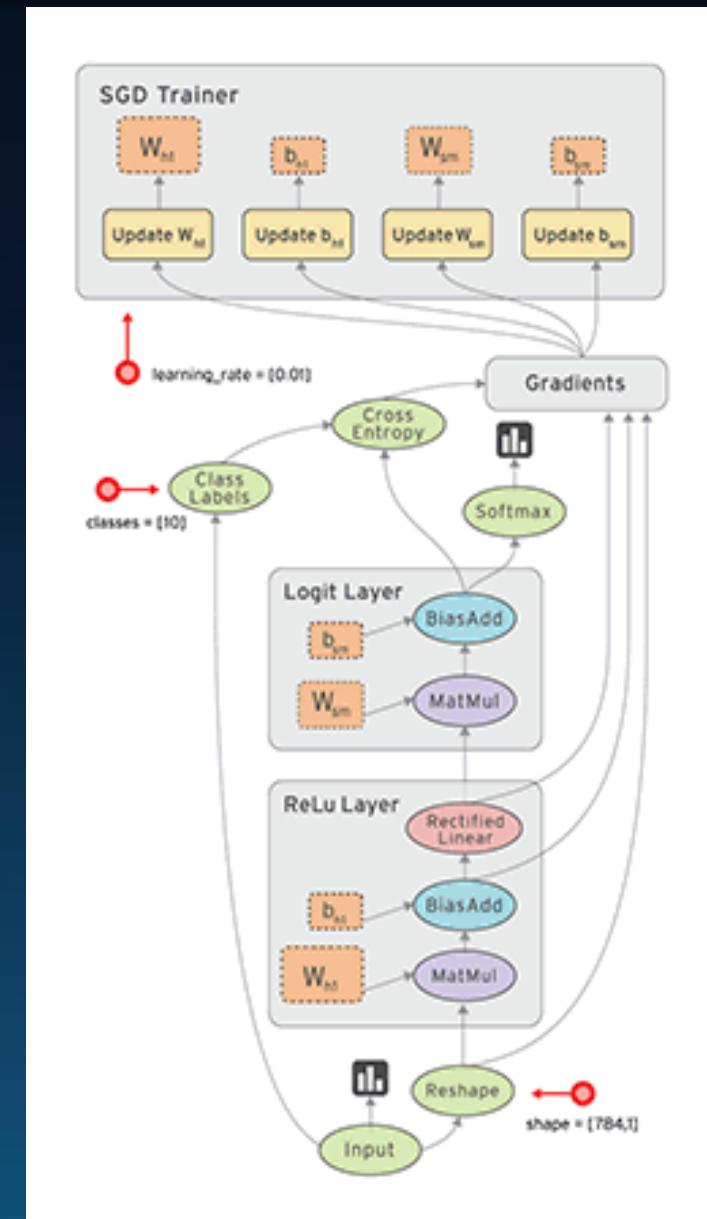
TensorFlow – Đứng từ quan điểm nghiên cứu viên

- Computational graph as **dataflow graph**
 - Construction phase:** assemble a graph to represent the design, and to train the model
 - Execution phase:** repeatedly execute a set of training ops in the graph



2-layer Neural Net

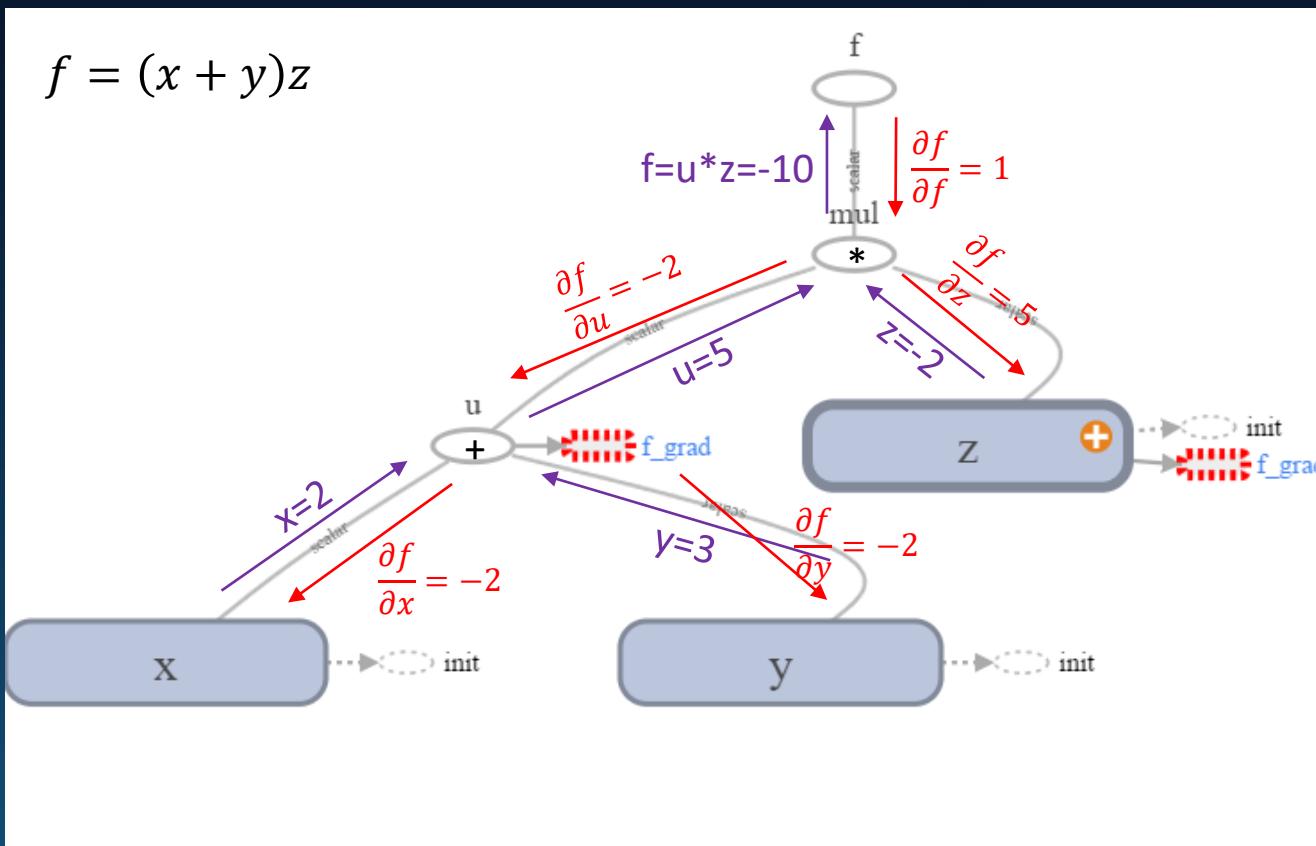
- Flexible, intuitive construction
- Launch the graph and run the training ops in a loop



Công nghệ TensorFlow

Một ví dụ tính toán đạo hàm tự động với TensorFlow

- The most attractive feature: Automatic Differentiation
- Automatically add ops to calculate symbolic gradients of variables w.r.t loss function, then apply these gradients with an optimization algorithm.

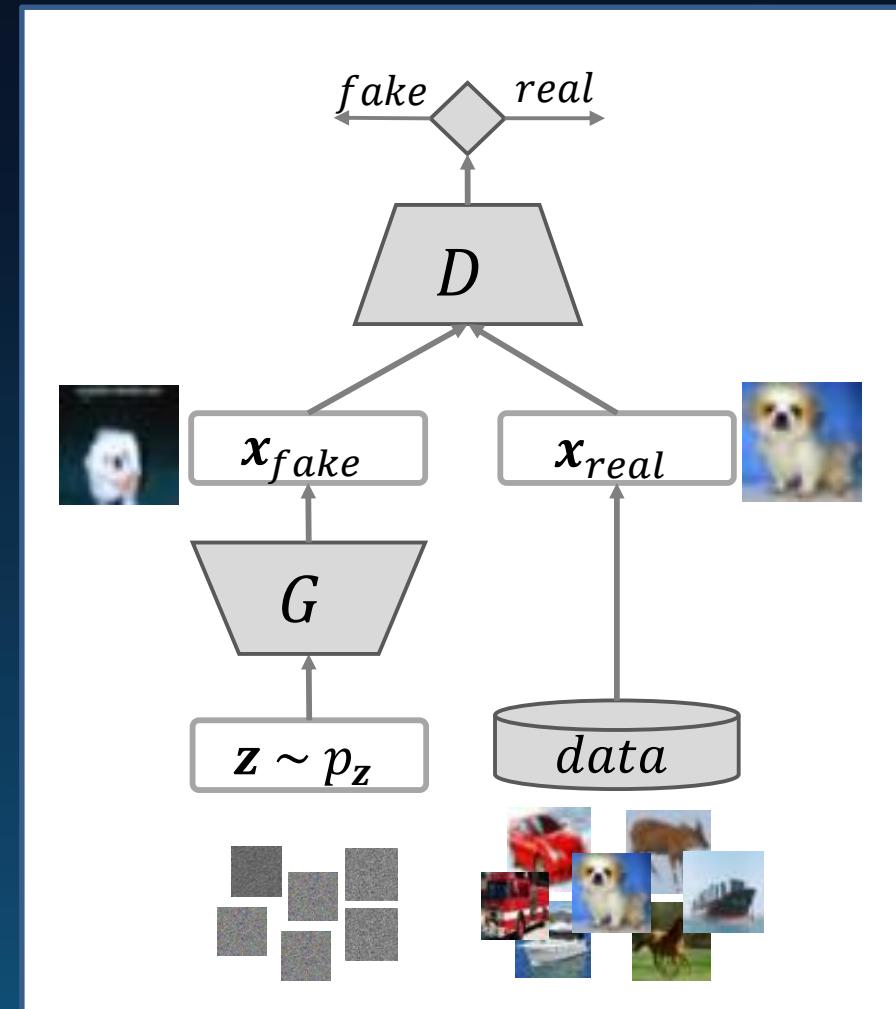
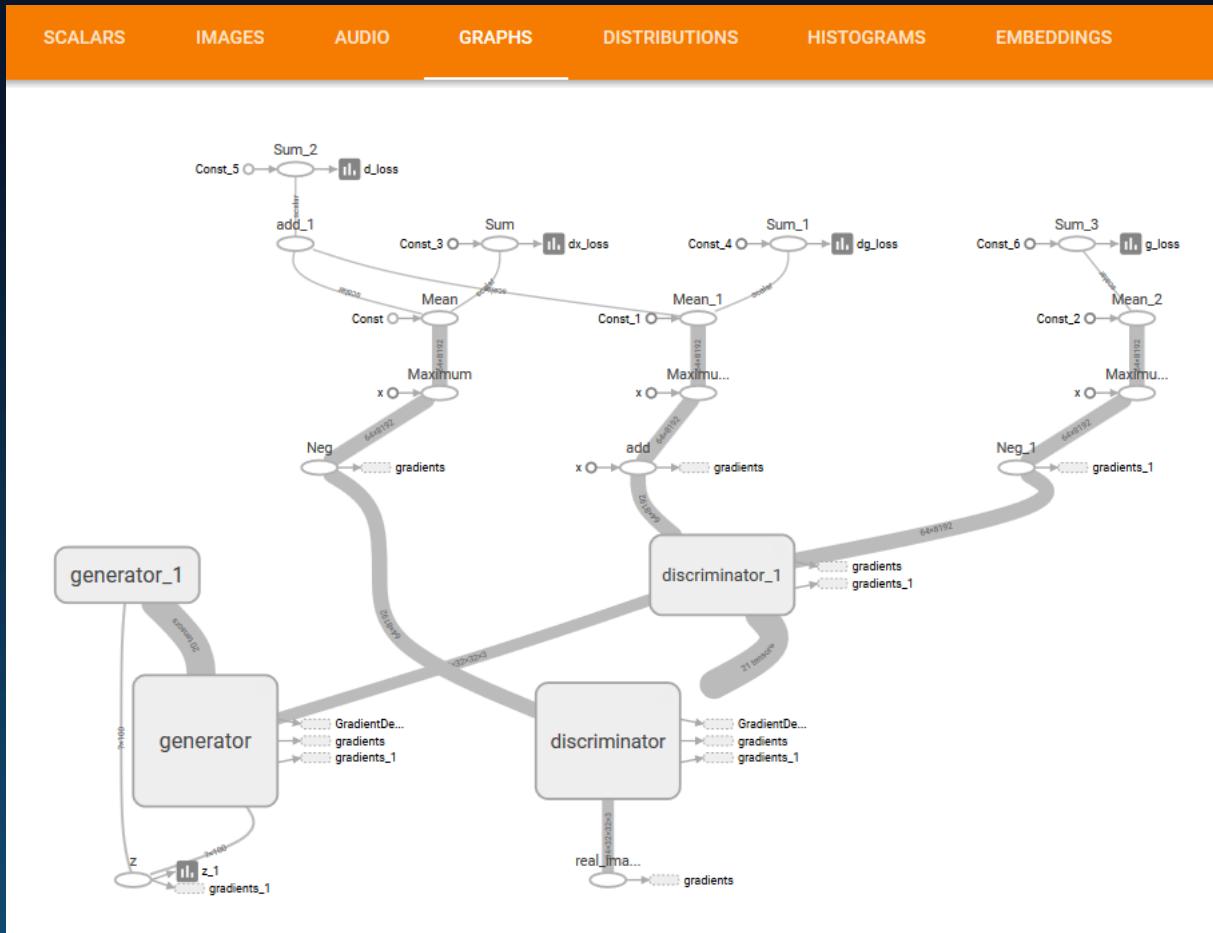


$$\begin{aligned}\frac{\partial f}{\partial f} &= 1 \\ \frac{\partial f}{\partial u} &= \frac{\partial f}{\partial f} \frac{\partial f}{\partial u} = 1 * 1 = 1 \\ \frac{\partial f}{\partial z} &= \frac{\partial f}{\partial f} \frac{\partial f}{\partial z} = 1 * z = 5 \\ \frac{\partial f}{\partial x} &= \frac{\partial f}{\partial u} \frac{\partial u}{\partial x} = 1 * 2 = 2 \\ \frac{\partial f}{\partial y} &= \frac{\partial f}{\partial u} \frac{\partial u}{\partial y} = 1 * 3 = 3\end{aligned}$$

Công nghệ TensorFlow

TensorFlow – Một số đồ thị tính toán cho mô hình nhiều tầng

- Computational graph of a Generative Adversarial Net (GAN)



Tổng kết

- Mô hình lớn cho dữ liệu lớn
 - Why traditional methods might fail on big data?
 - Three approaches to scale up big models
 - Scale up Gradient Descent Methods
 - Scale up Probabilistic Inference
 - Scale up Model Evaluation (bootstrap)
- Ba cách tiếp cận để phát triển mô hình lớn
 - Data augmentation
 - Stochastic Variational Inference for Graphical Models
 - Stochastic Gradient Descent and Online Learning
- Công nghệ nào cho mô hình lớn?
 - Hadoop, Spark, TensorFlow

Tài liệu tham khảo

Acknowledgement and References

- Một số hình ảnh minh họa được download từ images.google.com theo cài đặt mặc định của search engine này.
- Eric Xing and Qirong Ho, *A New Look at the System Algorithm and Theory Foundations of Distributed Machine Learning*, KDD Tutorial 2015.
- Michael Jordan, *On the Computational and Statistical Interface and “Big Data”*, ICML Keynote Speech, 2014.
- Feng Niu, Benjamin Recht, Christopher Ré, Stephen J. Wright, *HOGWILD!: A Lock-Free Approach to Parallelizing Stochastic Gradient Descent*, CoRR abs/1106.5730, 2011.
- Qirong Ho, James Cipar, Henggang Cui, Seunghak Lee, Jin Kyu Kim, Phillip B. Gibbons, Garth A. Gibson, Gregory R. Ganger, Eric P. Xing, *More Effective Distributed ML via a Stale Synchronous Parallel Parameter Server*, NIPS 2013: 1223-1231.
- Wei Dai, Abhimanu Kumar, Jinliang Wei, Qirong Ho, Garth A. Gibson, Eric P. Xing, *High-Performance Distributed ML at Scale through Parameter Server Consistency Models*. AAAI 2015: 79-87.
- Salomon Bochner, *Lectures on Fourier Integrals.(AM-42)*. Vol. 42. Princeton University Press, 2016.
- Zhuang Wang, Koby Crammer, Slobodan Vucetic, *Breaking the curse of kernelization: budgeted stochastic gradient descent for large-scale SVM training*, Journal of Machine Learning Research 13: 3103-3131, 2012.
- Jing Lu, Steven C. H. Hoi, Jialei Wang, Peilin Zhao, Zhi-Yong Liu, *Large Scale Online Kernel Learning*, Journal of Machine Learning Research (JMLR) 17(47) 1-43 2016.

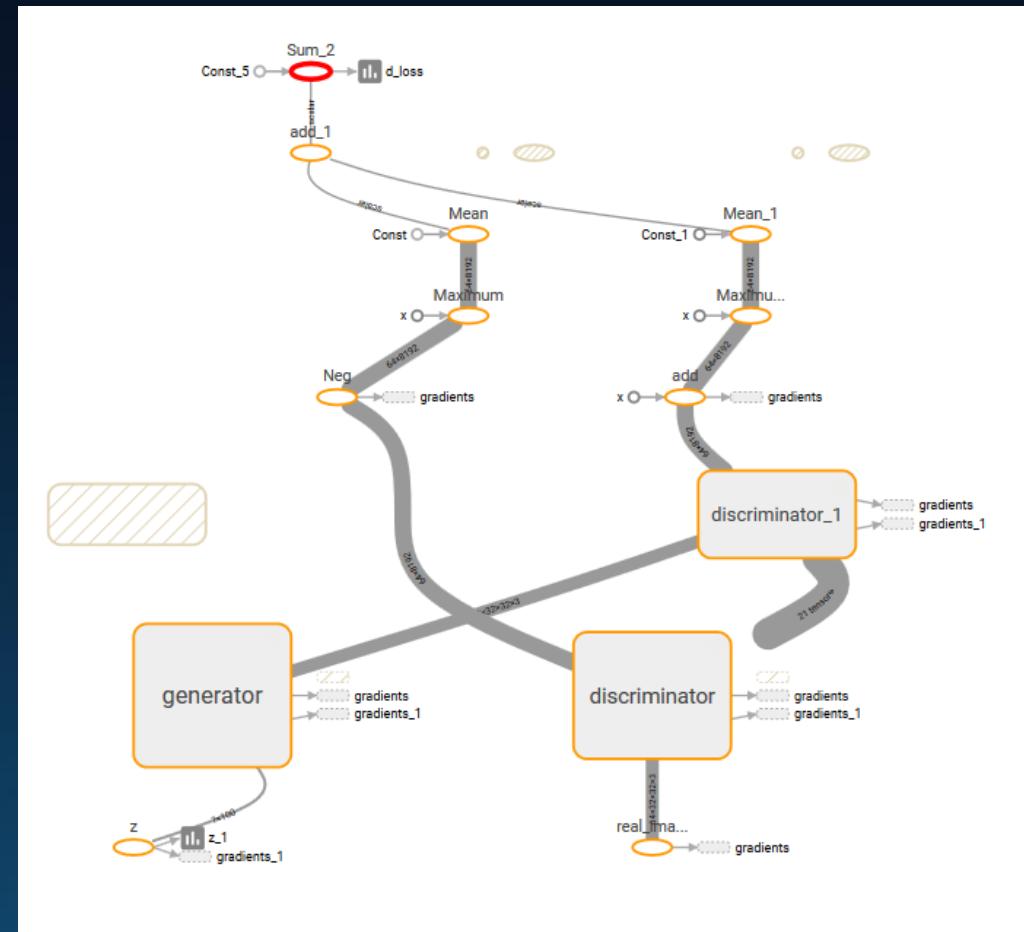
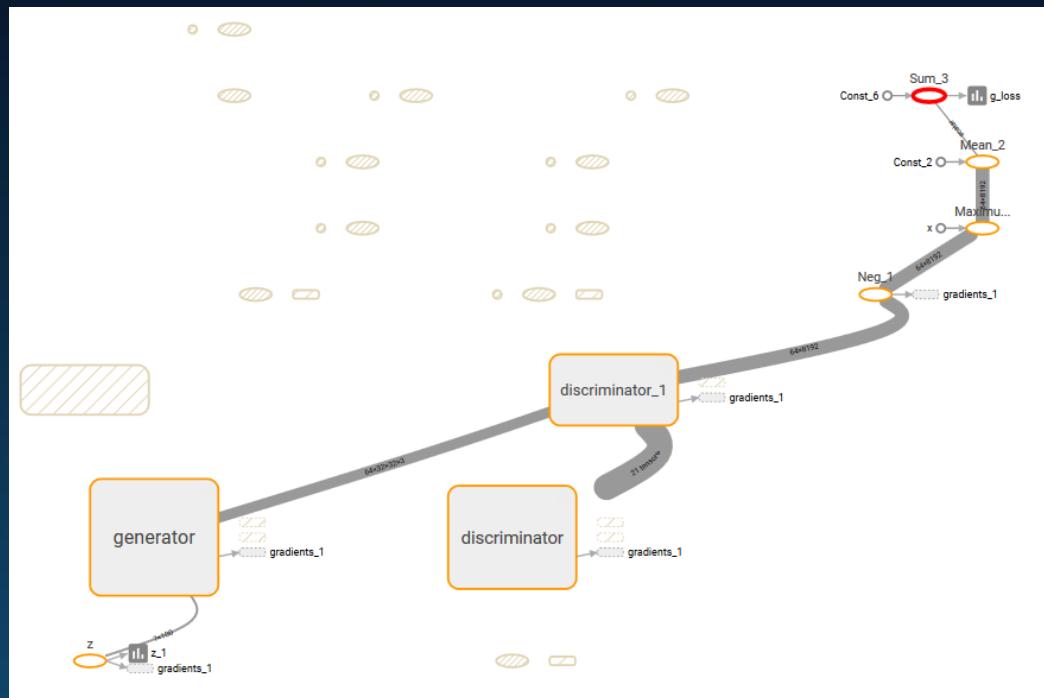
Xin cảm ơn các anh, chị đã lắng nghe!

Công nghệ TensorFlow

Debug với Tensorflow

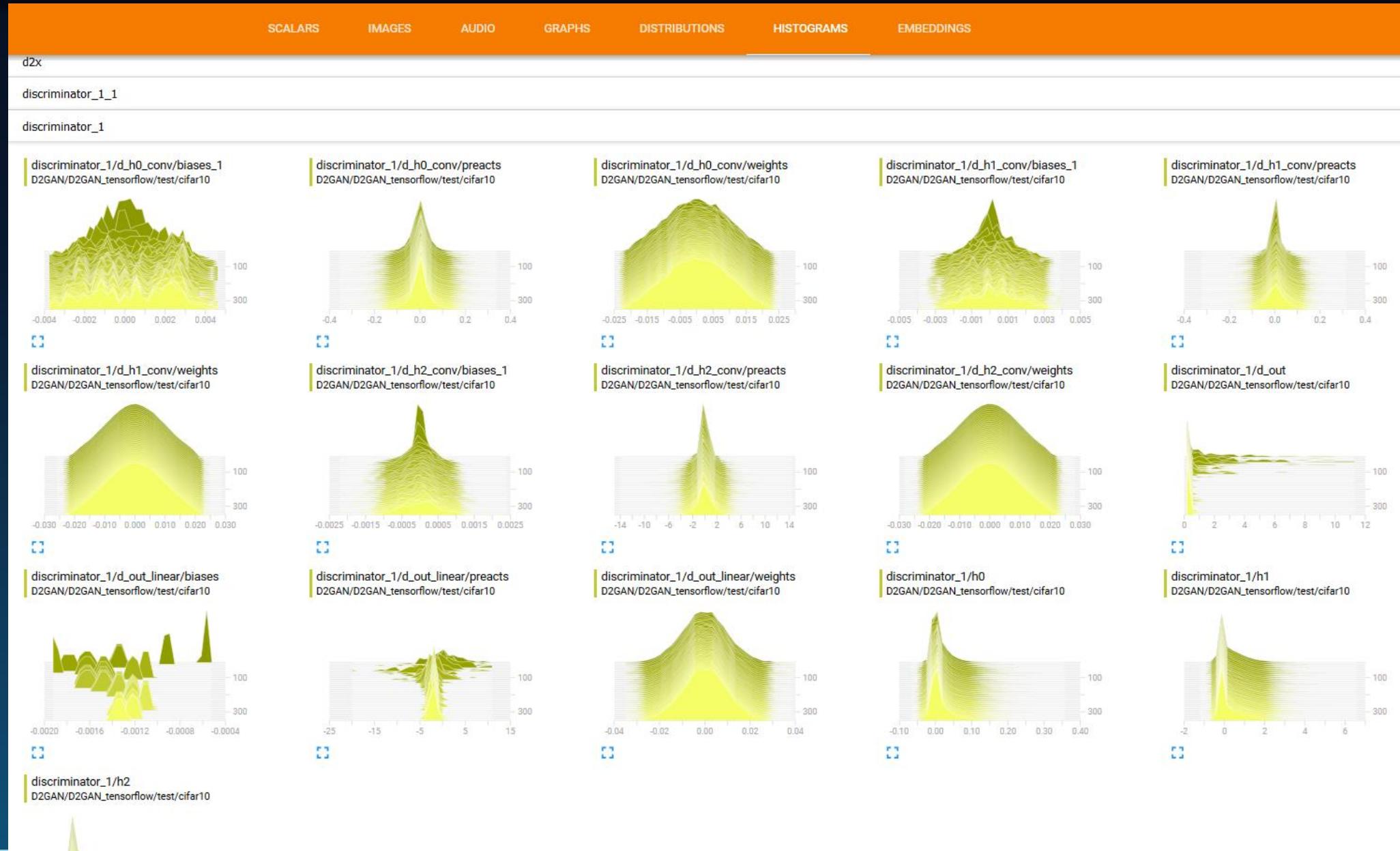
Computational Graph

- Trace connections and input



Công nghệ TensorFlow

Debug với Tensorflow



Debug với Tensorflow

