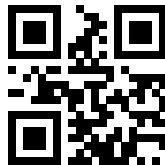# MassiveClicks

## A Massively-parallel Framework for Efficient Click Models Training

August 28, 2023

Skip Thijssen, Ana-Lucia Varbanescu, Pooya Khandel, Andrew Yates

**University of Amsterdam**

# What is a Click Model?

1. Users interact with **search engines**.
2. **Clicks** on search results show **relevance**.
3. Relevance assessments help **improve search engines**[1].

**Search**

| Search Query 🔍 |
| --- |
| **0**. <u>Search result</u> |
| **1**. <u>Search result</u> |
| **2**. <u>Search result</u> |
| **3**. <u>Search result</u> |
| **4**. <u>Search result</u> |
| **5**. <u>Search result</u> |
| **6**. <u>Search result</u> |
| **7**. <u>Search result</u> |
| **8**. <u>Search result</u> |
| **9**. <u>Search result</u> |

---

[1] Filip Radlinski et al. "How Does Clickthrough Data Reflect Retrieval Quality?"

# What is a Click Model?

**Click**

*Search Query* 🔍

**0**. <u>Search result</u>
**1**. <u>Search result</u>
**2**. <u>Search result</u>
**3**. <u>**Search result**</u>
**4**. <u>Search result</u>
**5**. <u>Search result</u>
**6**. <u>Search result</u>
**7**. <u>Search result</u>
**8**. <u>Search result</u>
**9**. <u>Search result</u>

1. Users interact with **search engines**.

2. **Clicks** on search results show **relevance**.
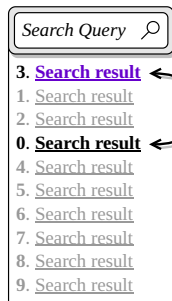
3. Relevance assessments help **improve search engines**[1].

---

[1]Filip Radlinski et al. "How Does Clickthrough Data Reflect Retrieval Quality?"

# What is a Click Model?

1. Users interact with **search engines**.
2. **Clicks** on search results show **relevance**.
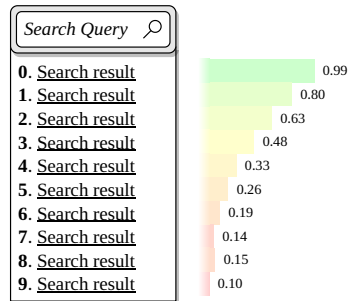3. Relevance assessments help **improve search engines**[1].

**Improve**

| Search Query 🔍 |

**3**. <u>**Search result**</u> ←
**1**. <u>Search result</u>
**2**. <u>Search result</u>
**0**. <u>**Search result**</u> ←
**4**. <u>Search result</u>
**5**. <u>Search result</u>
**6**. <u>Search result</u>
**7**. <u>Search result</u>
**8**. <u>Search result</u>
**9**. <u>Search result</u>

---

[1]Filip Radlinski et al. "How Does Clickthrough Data Reflect Retrieval Quality?"

# What is a Click Model?

A click model assigns a **click probability** to a search engine result page (SERP).

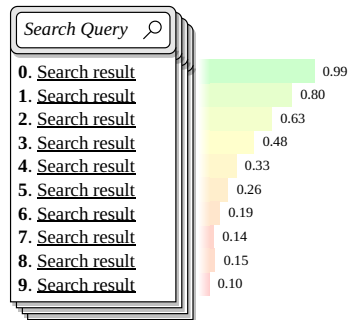# How is a Click Model built?

1. A search engine logs a user's clicks on a SERP inside a *click* log.
2. Click models use parameters to quantify relevance
3. EM-based models (for example) are used to train these parameters using data from click logs.

Click logs from real-world search engines can be very large[2] making training expensive.



| *Search Query* 🔍 | |
|---|---|
| **0**. <u>Search result</u> | 0.99 |
| **1**. <u>Search result</u> | 0.80 |
| **2**. <u>Search result</u> | 0.63 |
| **3**. <u>Search result</u> | 0.48 |
| **4**. <u>Search result</u> | 0.33 |
| **5**. <u>Search result</u> | 0.26 |
| **6**. <u>Search result</u> | 0.19 |
| **7**. <u>Search result</u> | 0.14 |
| **8**. <u>Search result</u> | 0.15 |
| **9**. <u>Search result</u> | 0.10 |

---

[2] https://www.internetlivestats.com/google-search-statistics/

# Limitations of Existing Solutions

— **Training** EM-based models is **challenging due to the size** of click logs.
— **Existing tools** like *PyClick* are sequential and slow.
— **_ParClick_**, though parallel, is limited to single-node multi-core systems.

# MassiveClicks

*A scalable multi-GPU solution to parallelize generic EM-based training algorithms for click models.*

**Why MassiveClicks?**

— First multi-GPU, distributed click model training framework.

— Efficient GPU kernels and data-partitioning.

— Outperforms ParClick on GPUs/multi-node.

**Requirements**

— Scalable EM-based training.

— Efficient, multi-GPU, multi-node distributed processing.

— General, adaptable framework.

*A scalable multi-GPU solution to parallelize generic EM-based training algorithms for click models.*

*Multi-node / Multi-GPU*



Node 1   Node 2   Node 3

# EM-based Click Models

Click models supported by MassiveClicks:

*Computational difficulty:*

- Position-based Model (PBM)                    8 %
- Click Chain Model (CCM)                       88 %
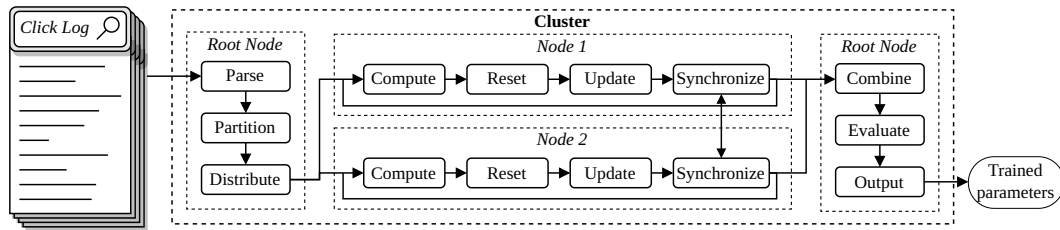- Dynamic Bayesian Network model (DBN)         100 %
- User Browsing Model (UBM)                      10 %

# EM-based Click Models

Click models supported by MassiveClicks:

*Computational
difficulty:*

- **Position-based Model (PBM)** **8 %**
- **Click Chain Model (CCM)** **88 %**
- Dynamic Bayesian Network model (DBN) 100 %
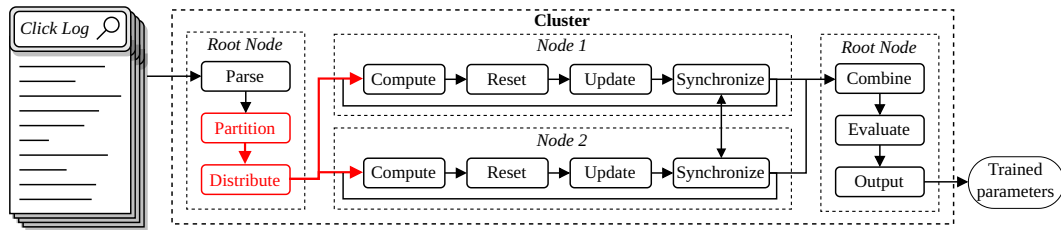- User Browsing Model (UBM) 10 %

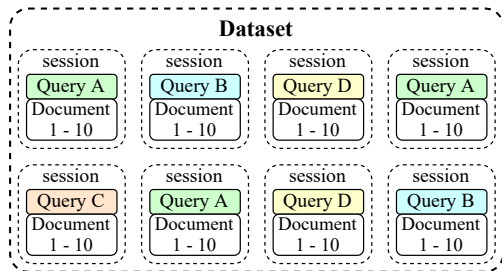1. Read click log.

2. Preprocess input.

3. Estimate parameters.

4. Evaluate and output results.

# Data Distribution

Partitioning and distribution

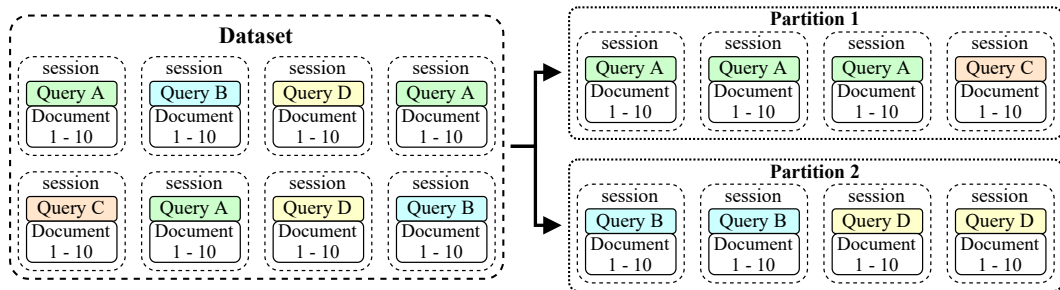# Data Distribution

Constrain data distribution to **reduce communication** during parameter estimation. Sessions are *grouped by query*.

# Data Distribution

Constrain data distribution to **reduce communication** during parameter estimation. Sessions are *grouped by query*.

# Data Distribution

Nodes *parse sessions independently* from root, ensuring dataset size isn't limited by the root node's memory.

**Difficulties:**

– Similar queries grouped on same node.

– Query IDs on other nodes unknown.

– IDs are inconsistent due to gaps.

– Variable query count per node.

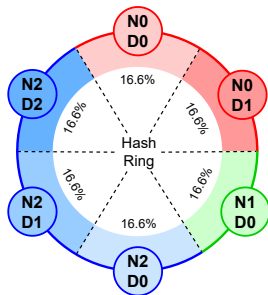– Differing memory available per node.

**Desirable Characteristics:**

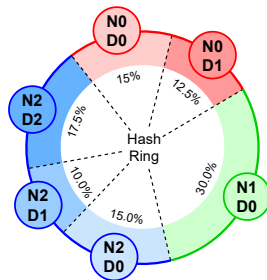– Minimize inter-node communication.

# Data Distribution

Decide session distribution across nodes by adjusting *node ranges* on a *hash ring* based on some property, i.e., number of CUDA cores.
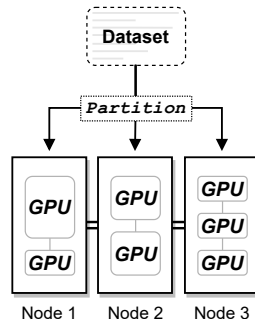


**Uniform Distribution**

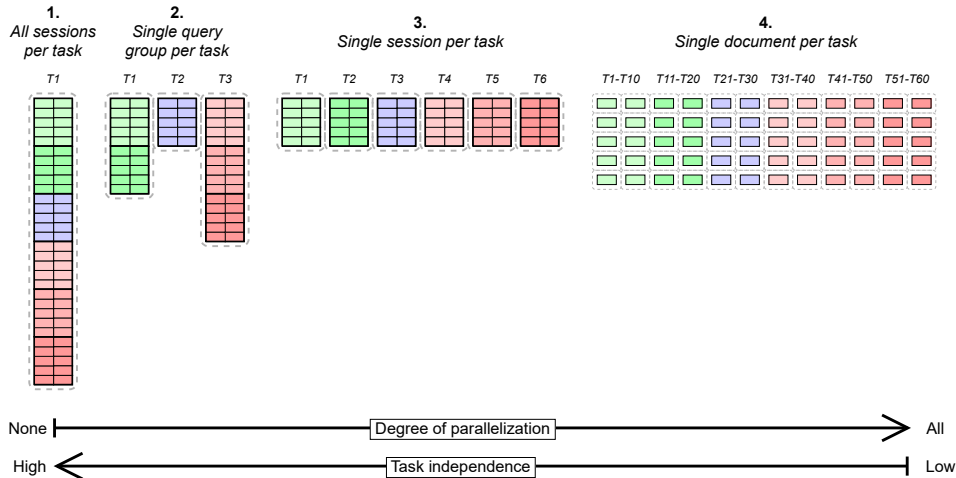**Property-Based Distribution**

# Data Distribution

Five partitioning policies:

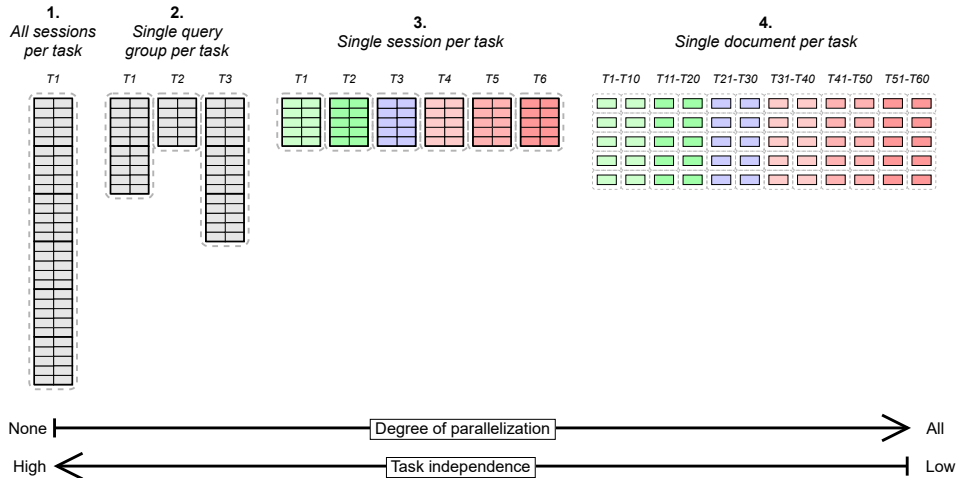- Round-robin

- Maximum-Utilization

- Proportional to:
    - Available Memory
    - CUDA-core Count
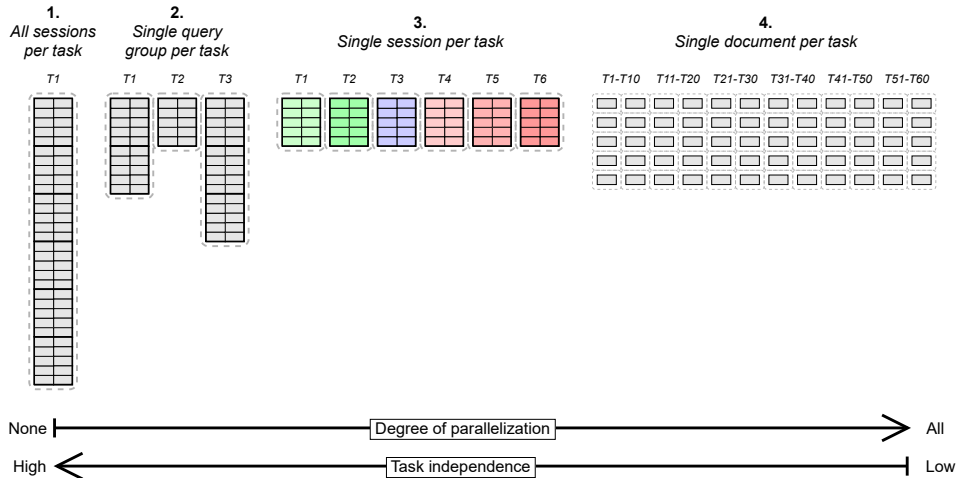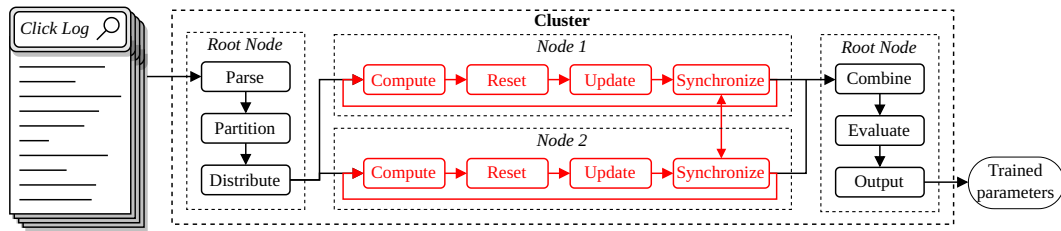    - Theoretical Peak Performance

# Session Distribution



**1.**
*All sessions per task*

T1

**2.**
*Single query group per task*

T1   T2   T3

**3.**
*Single session per task*

T1   T2   T3   T4   T5   T6

**4.**
*Single document per task*

T1-T10   T11-T20   T21-T30   T31-T40   T41-T50   T51-T60

None |———————————————— Degree of parallelization ————————————————▶ All

High ◀———————————————— Task independence ————————————————| Low

# Session Distribution



**1.**
*All sessions per task*

**2.**
*Single query group per task*

**3.**
*Single session per task*

**4.**
*Single document per task*

T1 — T1 T2 T3 — T1 T2 T3 T4 T5 T6 — T1-T10 T11-T20 T21-T30 T31-T40 T41-T50 T51-T60

None ├──────── Degree of parallelization ────────► All

High ◄──────── Task independence ────────┤ Low

# Session Distribution

**1.**
*All sessions per task*

**2.**
*Single query group per task*

**3.**
*Single session per task*

**4.**
*Single document per task*

*T1*

*T1   T2   T3*

*T1   T2   T3   T4   T5   T6*

*T1-T10   T11-T20   T21-T30   T31-T40   T41-T50   T51-T60*

None |———————————— Degree of parallelization ————————————▶ All

High ◀———————————— Task independence ————————————| Low

# Parameter Estimation

Iterative parameter estimation

# Parameter Estimation

EM-based parameter estimation is divided into four phases.
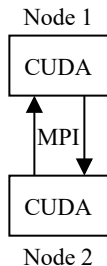
# Parameter Estimation

EM-based parameter estimation is divided into four phases.
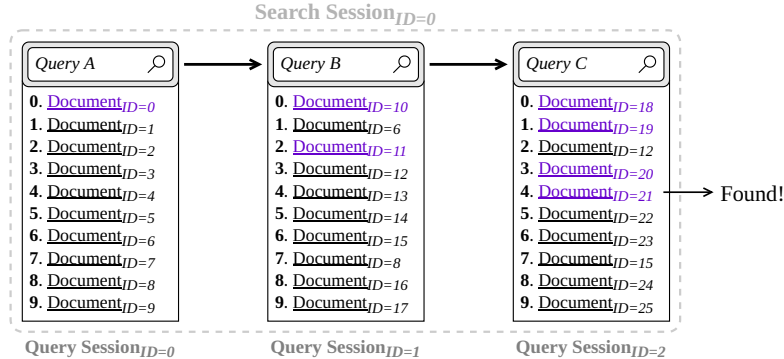
# Multi-GPU/Multi-node Architectures

Programming models:

- **C++** framework.
- **CUDA** for (multi) GPU execution.
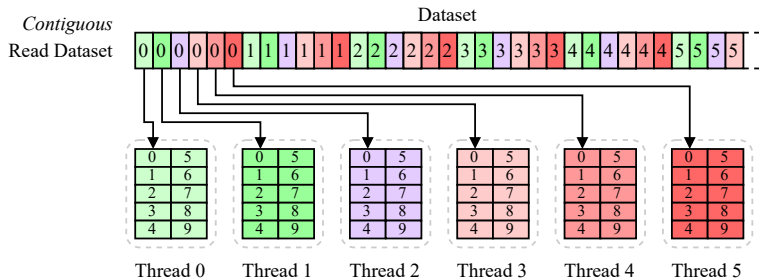- **MPI** for multi-node communication.



Node 1

CUDA
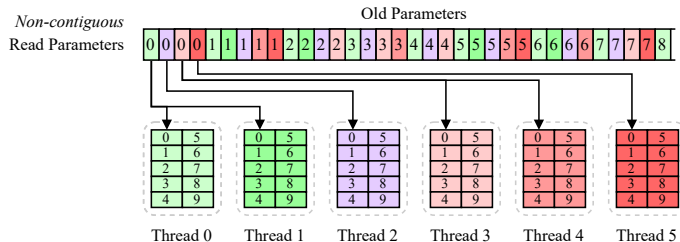
MPI

CUDA

Node 2

# Click Log Layout

**Search Session**$_{ID=0}$

| Query A $\mathcal{Q}$ | Query B $\mathcal{Q}$ | Query C $\mathcal{Q}$ |
|---|---|---|
| **0**. Document$_{ID=0}$ | **0**. Document$_{ID=10}$ | **0**. Document$_{ID=18}$ |
| **1**. Document$_{ID=1}$ | **1**. Document$_{ID=6}$ | **1**. Document$_{ID=19}$ |
| **2**. Document$_{ID=2}$ | **2**. Document$_{ID=11}$ | **2**. Document$_{ID=12}$ |
| **3**. Document$_{ID=3}$ | **3**. Document$_{ID=12}$ | **3**. Document$_{ID=20}$ |
| **4**. Document$_{ID=4}$ | **4**. Document$_{ID=13}$ | **4**. Document$_{ID=21}$ → Found! |
| **5**. Document$_{ID=5}$ | **5**. Document$_{ID=14}$ | **5**. Document$_{ID=22}$ |
| **6**. Document$_{ID=6}$ | **6**. Document$_{ID=15}$ | **6**. Document$_{ID=23}$ |
| **7**. Document$_{ID=7}$ | **7**. Document$_{ID=8}$ | **7**. Document$_{ID=15}$ |
| **8**. Document$_{ID=8}$ | **8**. Document$_{ID=16}$ | **8**. Document$_{ID=24}$ |
| **9**. Document$_{ID=9}$ | **9**. Document$_{ID=17}$ | **9**. Document$_{ID=25}$ |
| **Query Session**$_{ID=0}$ | **Query Session**$_{ID=1}$ | **Query Session**$_{ID=2}$ |

| Session $ID=0$ | Time passed | Query | Query $ID=0$ | Region ID | Document $ID=0$ | Document $ID=1$ | Document $ID=2$ | Document $ID=3$ | Document $ID=4$ | Document $ID=5$ | Document $ID=6$ | Document $ID=7$ | Document $ID=8$ | Document $ID=9$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Session $ID=0$ | Time passed | Click | Document $ID=0$ | | | | | | | | | | | |

# Allocation in Memory

Dataset placement in global GPU memory.
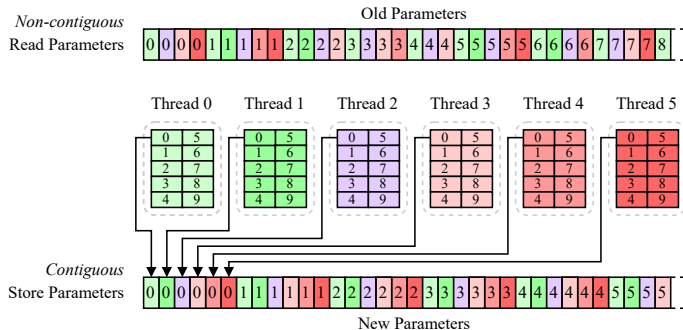
# Memory Access Pattern

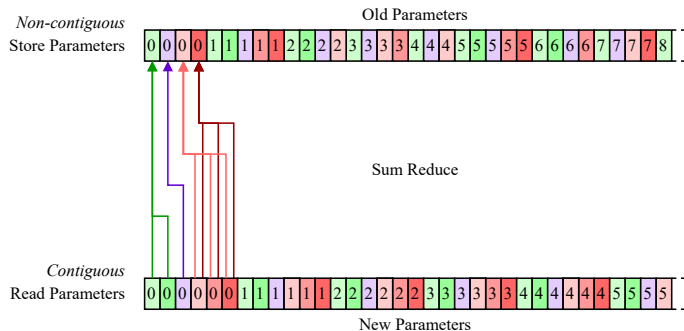**Compute** — Threads *read* parameters from *previous* iteration.

# Memory Access Pattern

**Compute** — Threads *write* parameters from *current* iteration.

# Memory Access Pattern

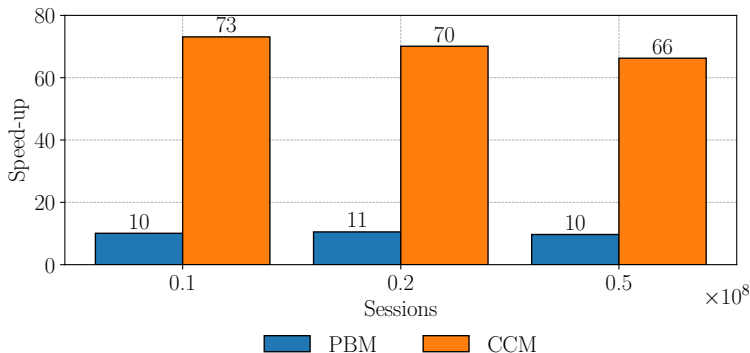**Update** — Threads *write* parameters to *previous* iteration.

# Experimental Setup

Measure performance for (up to 14) multi-GPU/multi-node configurations using **NVIDIA RTX A4000 GPUs** and 24-core **AMD EPYC 7402P CPUs** and subsets of the Yandex dataset of varying sizes ($Dn$ = dataset of $n$ million sessions).

Report:

- **Speed-up** vs. ParClick, the **only alternative** for parallel click model training.
- **Scalability** for multi-GPU/multi-node setups and click logs.
- **Usability** for real-world problems.
- **Kernel performance** using a **roofline model**.

# Single GPU Speed-up

**Speed-up** of MassiveClicks for **PBM** and **CCM** on a single **NVIDIA RTX A4000** GPU.

# Multi-GPU Speed-up



**Speed-up** for **PBM** (left) and **CCM** (right) for different datasets and **NVIDIA RTX A4000 GPUs** compared to ParClick on an **AMD EPYC 7402P CPU** with **48 threads**.
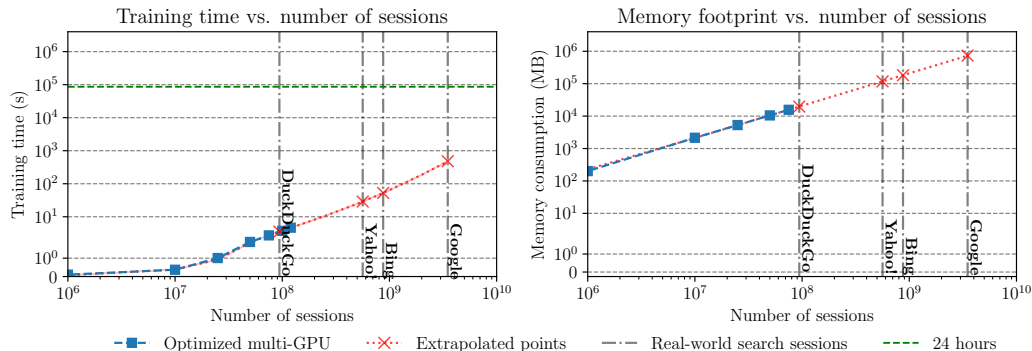
**Scalability** for **PBM** (left) and **CCM** (right) for different datasets and **NVIDIA RTX A4000 GPUs** compared to ParClick on an **AMD EPYC 7402P CPU** with **48 threads**.

**Training time** of MassiveClicks for **PBM** (left) and **CCM** (right) for **D10** with different number of **NVIDIA RTX A4000** GPUs per node.

Real-world applications:

  – 'Search engine'-size click logs.

  – Energy consumption.

**Training time** (left) and **memory footprint** (right) of MassiveClicks of **PBM** for real-world datasets using 14 **NVIDIA RTX A4000** GPUs.
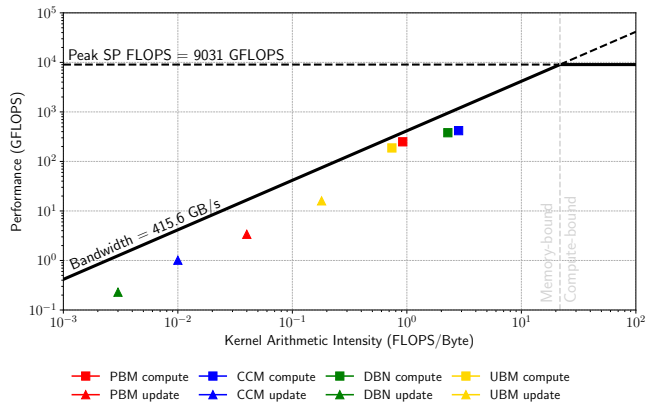
# Projection



**Training time** (left) and **memory footprint** (right) of MassiveClicks of **CCM** for real-world datasets using 14 **NVIDIA RTX A4000** GPUs.

# Energy Consumption



Expected **energy consumption** of MassiveClicks for **PBM** and **CCM** compared to ParClick.
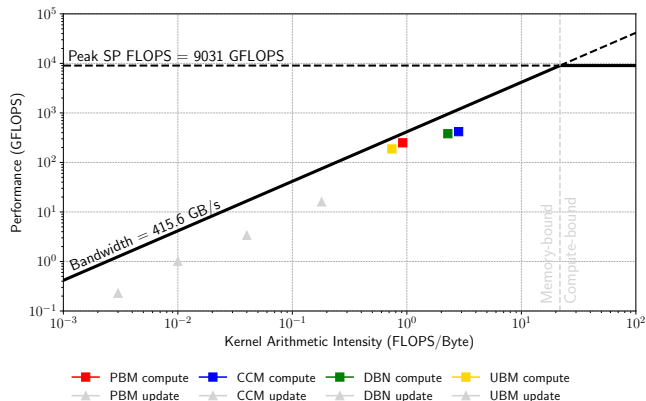
# Kernel Performance

**Roofline model** of MassiveClicks for **PBM**, **CCM**, **DBN**, and **UBM** on a single **NVIDIA RTX A4000** GPU.
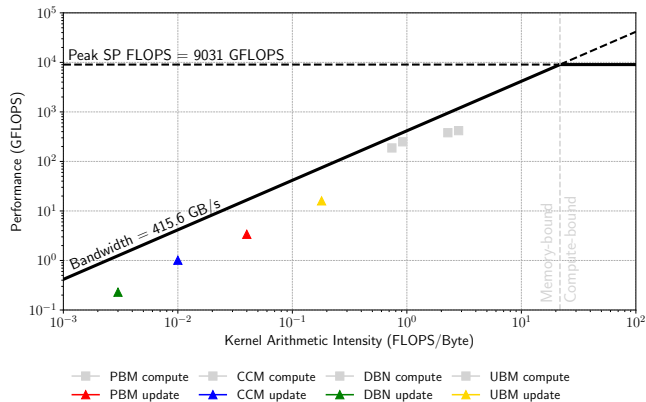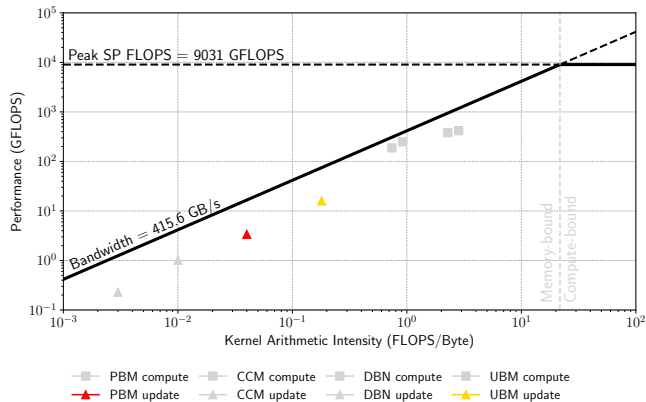
**Roofline model** of MassiveClicks for **PBM**, **CCM**, **DBN**, and **UBM** on a single **NVIDIA RTX A4000** GPU.

**Roofline model** of MassiveClicks for **PBM**, **CCM**, **DBN**, and **UBM** on a single **NVIDIA RTX A4000** GPU.

**Roofline model** of MassiveClicks for **PBM**, **CCM**, **DBN**, and **UBM** on a single **NVIDIA RTX A4000** GPU.

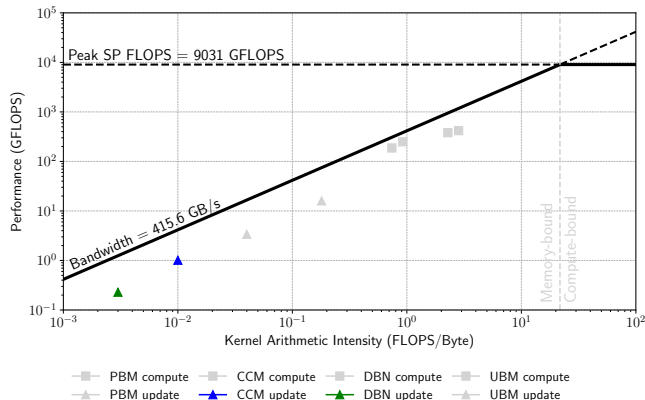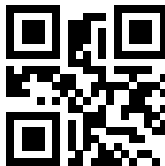**Roofline model** of MassiveClicks for **PBM**, **CCM**, **DBN**, and **UBM** on a single **NVIDIA RTX A4000** GPU.

# Conclusion and Future Directions

- **In Summary:** We introduced MassiveClicks, a tool for training EM-based click models on heterogeneous multi-GPU/multi-node setups, offering a significant performance improvement over existing solutions.

- **Future Work:**
    - Exploring hybrid training methods.
    - Transition to HIP for broader compatibility.

- **Repository:** Find the code and documentation for MassiveClicks at *github.com/skip-th/MassiveClicks* or the QR-code.

# Hybrid training

### Automatic

  – Compute optimal data distribution.

### Manual

  – User chooses desired distribution.

# Compatibility

The heterogeneity only extends to CUDA GPUs.

– AMD GPUs are not supported.

– Currently converting codebase to HIP using HIPify[1].

---

[1]https://github.com/ROCm-Developer-Tools/HIPIFY

# Position-based Model (PBM)

$$P(C_d = 1) = P(E_d = 1) \cdot P(A_d = 1)$$
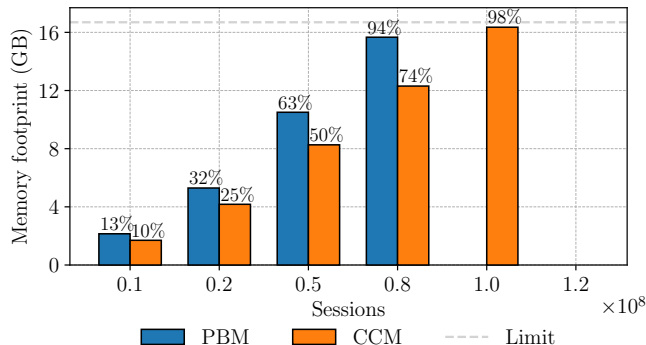
A user **clicks** a search result if, and only if, they **examined** the result and were **attracted** to it.

$\gamma_r$ – rank-dependent *examination* parameters.

$\alpha_{qd}$ – query-dependent *attractiveness* parameter.

# Click Chain Model (CCM)

$$P(C_d = 1) = P(E_d = 1) \cdot P(A_d = 1)$$

A user **clicks** a search result if, and only if, they **examined** the result and were **attracted** to it.

$\tau_1, \tau_2, \tau_3$ — document-dependent *continuation* parameters.

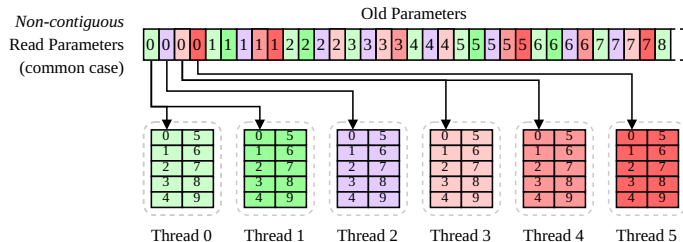$\alpha_{q_r d}$ — query-dependent *relevance* parameter.

# Memory Footprint



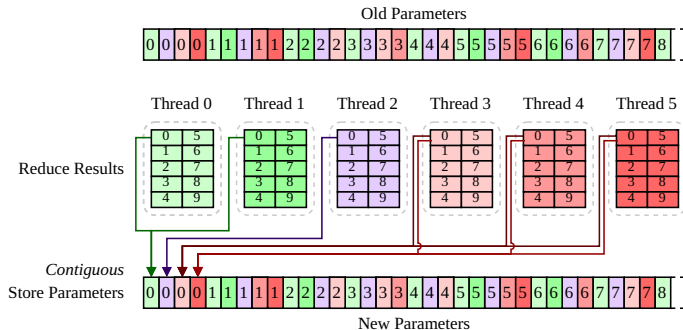**Memory footprint** of MassiveClicks for **PBM** and **CCM** on a single **NVIDIA RTX A4000** GPU.

**Compute** — Threads *read* parameters from *previous* iteration.

# New Memory Access Pattern

**Update** — Threads *reduce* parameters to *current* iteration.

# New Memory Access Pattern

**Update** — Threads *swap* old and new parameters.