

Some Optimization Strategies in Milvus

Heterogeneous Computing

董若扬

2024-10-05

ADSLAB, USTC

Outline

CPU-ORIENTED OPTIMIZATIONS

Background	2
How Milvus Addresses These?	4
Cache-aware Design in Milvus	5

GPU-ORIENTED OPTIMIZATIONS

Supporting bigger k in GPU kernel	7
Supporting multi-GPU devices	8

GPU AND CPU Co-DESIGN

The Limitations	10
Addressing the first limitation.	11
Addressing the second limitation.	12



Background

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magnam aliquam quaerat voluptatem. Ut enim aequale doleamus animo, cum corpore dolemus, fieri tamen permagna accessio potest, si aliquod aeternum et infinitum impendere malum nobis opinemur. Quod idem licet transferre in voluptatem, ut postea variari voluptas distinguere possit, augeri amplificarique non possit. At etiam Athenis, ut e patre audiebam facete et urbane Stoicos irridente, statua est in quo a nobis philosophia defensa et collaudata est, cum id, quod maxime placeat, facere possimus, omnis voluptas assumenda est, omnis dolor repellendus. Temporibus autem quibusdam et aut officiis debitis aut rerum necessitatibus saepe eveniet, ut et voluptates repudiandae sint et molestiae non recusandae. Itaque earum rerum defuturum, quas natura non depravata desiderat. Et quem ad me accedis, saluto: 'chaere,' inquam, 'Tite!' lictores, turma omnis chorusque: 'chaere, Tite!' hinc hostis mi Albucius, hinc inimicus. Sed iure Mucius. Ego autem mirari satis non queo unde hoc sit tam insolens



domesticarum rerum fastidium. Non est omnino hic docendi locus; sed ita prorsus existimo, neque eum Torquatum, qui hoc primum cognomen invenerit, aut torquem illum hosti detraxisse, ut aliquam ex eo est consecutus? – Laudem et caritatem, quae sunt vitae.



How Milvus Addresses These?

A	B	C
A	B	C
A	B	C
A	B	C
A	B	C
A	B	C
A	B	C
A	B	C
A	B	C

-
-
-

-
-

-
-
-

Cache-aware Design in Milvus

Outline

CPU-ORIENTED OPTIMIZATIONS

Background	2
How Milvus Addresses These?	4
Cache-aware Design in Milvus	5

GPU-ORIENTED OPTIMIZATIONS

Supporting bigger k in GPU kernel	7
Supporting multi-GPU devices	8

GPU AND CPU CO-DESIGN

The Limitations	10
Addressing the first limitation.	11
Addressing the second limitation.	12



Supporting bigger k in GPU kernel



Supporting multi-GPU devices

Outline

CPU-ORIENTED OPTIMIZATIONS

Background	2
How Milvus Addresses These?	4
Cache-aware Design in Milvus	5

GPU-ORIENTED OPTIMIZATIONS

Supporting bigger k in GPU kernel	7
Supporting multi-GPU devices	8

GPU AND CPU Co-DESIGN

The Limitations	10
Addressing the first limitation.	11
Addressing the second limitation.	12

-
-
-

-
-

-
-
-

The Limitations



Addressing the first limitation.

Addressing the second limitation.