# Some Optimization Strategies in Milvus

## Heterogeneous Computing

董若扬

2024-10-02

ADSLAB, USTC

# Outline

# Outline

# 1.1 Background

# 1.2 How Milvus Addresses These?

# 1.3 Cache-aware Design in Milvus

# Outline

# 2.1 Supporting bigger k in GPU kernel

## 2.2 Supporting multi-GPU devices

# Outline

# 3.1 The Limitations

## 3.2 Addressing the first limitation.

# 3.3 Addressing the second limitation.