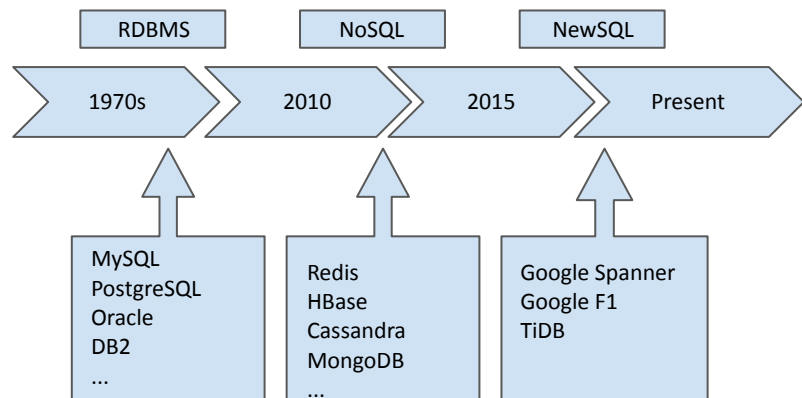# Go in TiDB

Yao Wei | PingCAP

# About me

- Yao Wei (姚维)

- TiDB Kernel Expert, General Manager of South Region, China

- 360 Infra team / Alibaba-UC / PingCAP

- Atlas/MySQL-Sniffer
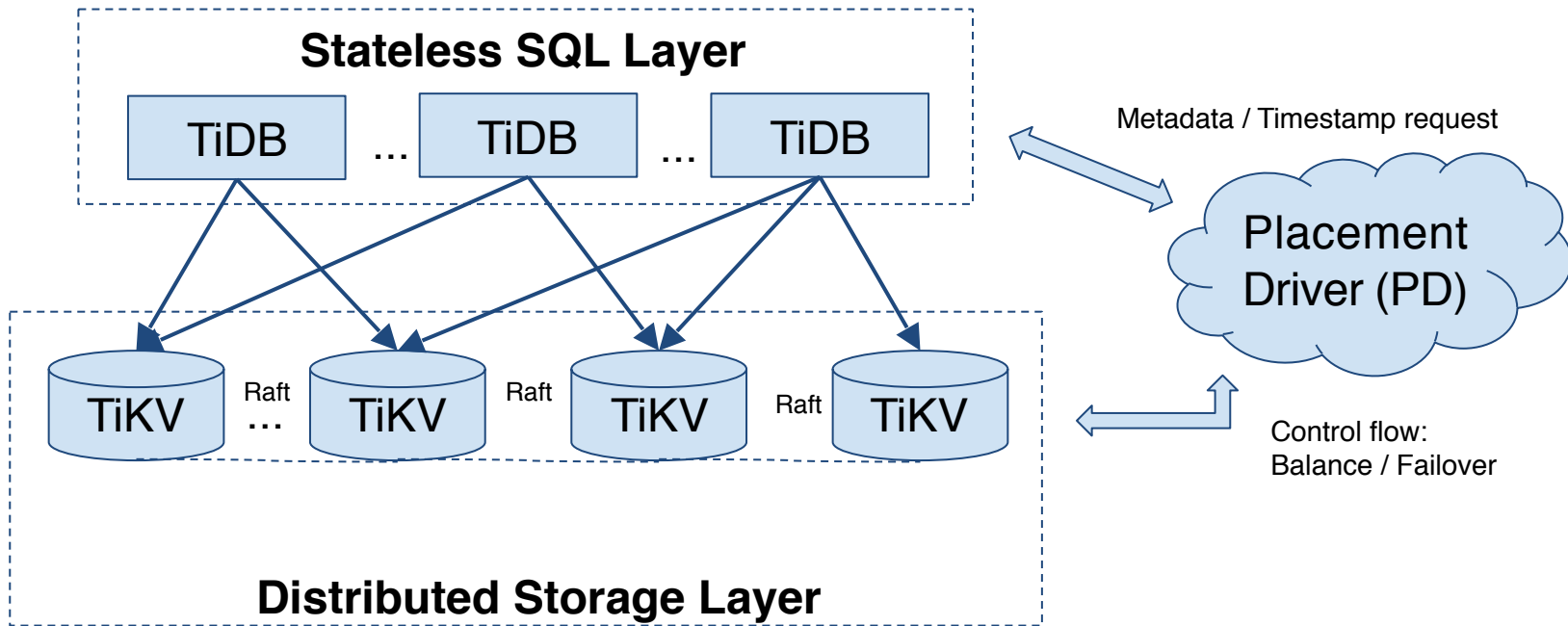
- Infrastructure software engineer

# Why a new database?

# Brief History

- Standalone RDBMS
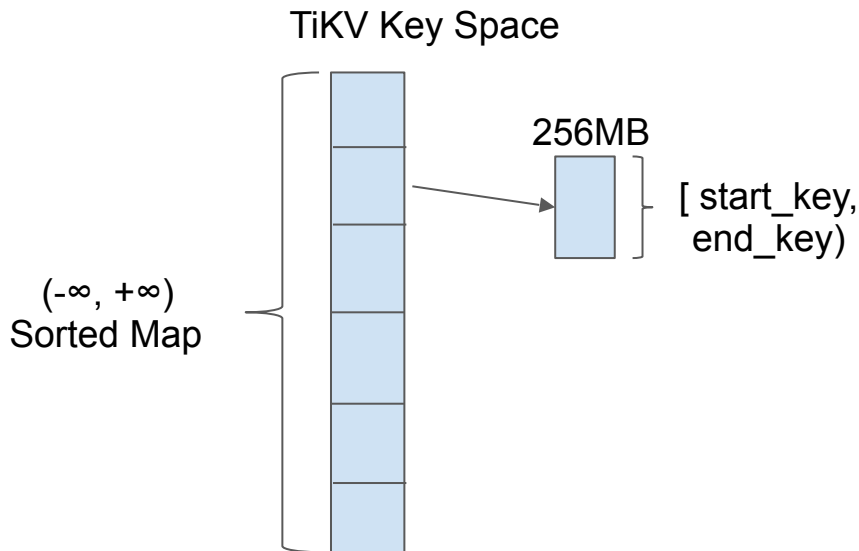- NoSQL
- Middleware & Proxy
- NewSQL

| RDBMS | NoSQL | NewSQL | |
|---|---|---|---|
| 1970s | 2010 | 2015 | Present |

MySQL
PostgreSQL
Oracle
DB2
...

Redis
HBase
Cassandra
MongoDB
...

Google Spanner
Google F1
TiDB

# Architecture



**Stateless SQL Layer**

TiDB … TiDB … TiDB

Metadata / Timestamp request

**Placement Driver (PD)**

**Distributed Storage Layer**

TiKV … Raft TiKV Raft TiKV Raft TiKV

Control flow:
Balance / Failover
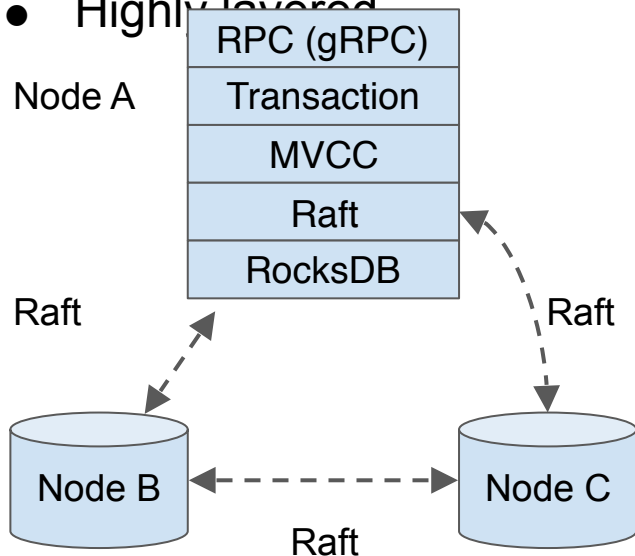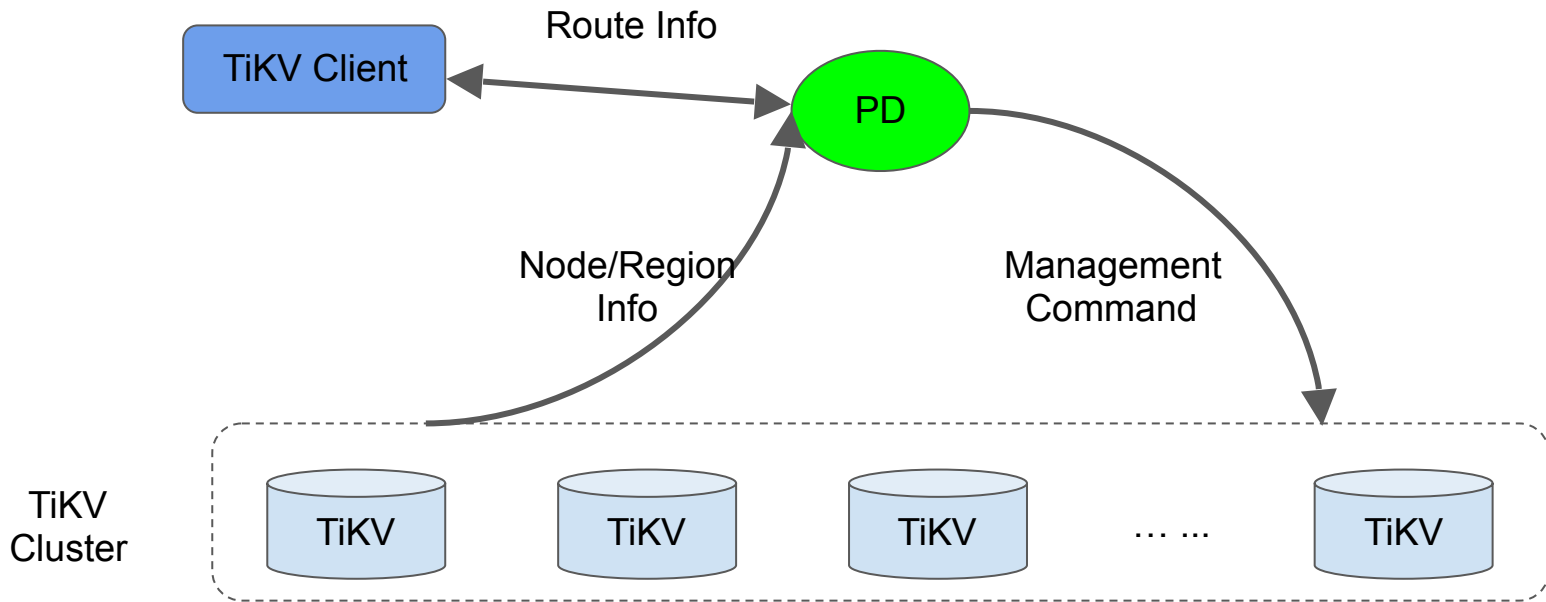
# TiKV - Overview

- Region: a set of continuous key-value pairs

- Data is organized/stored/replicated by Regions

- Highly layered

Node A

| RPC (gRPC) |
| Transaction |
| MVCC |
| Raft |
| RocksDB |

Raft

Raft

Node B

Node C

Raft

TiKV Key Space
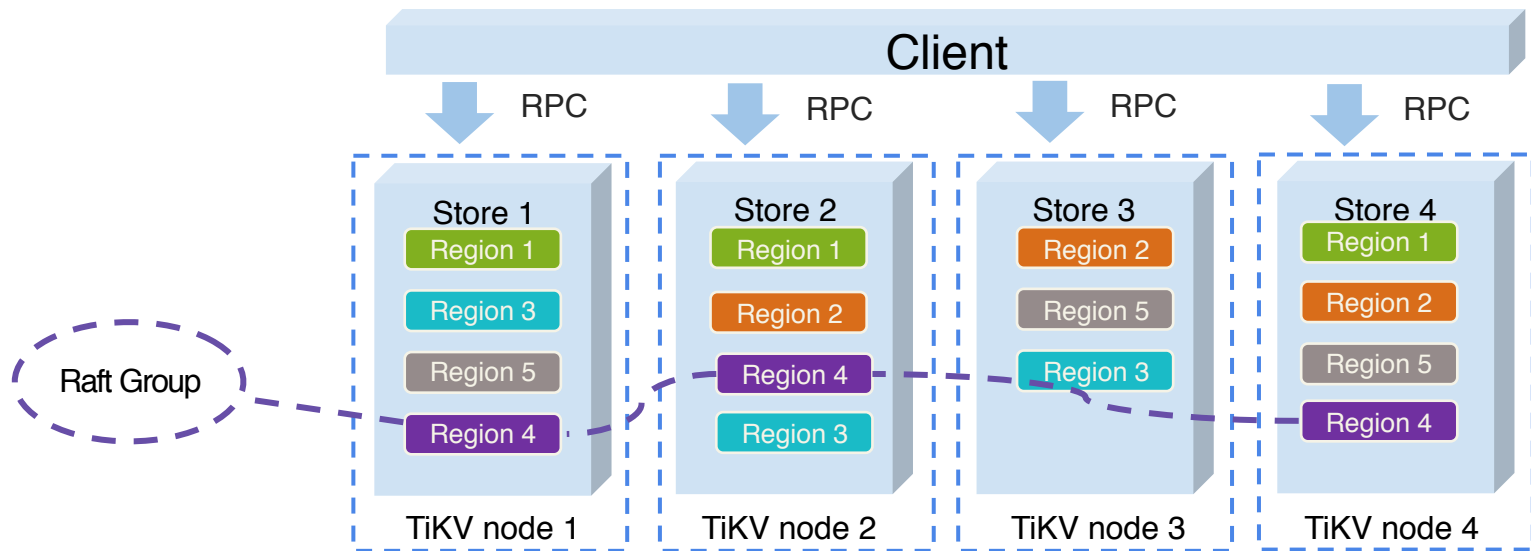
256MB

[ start_key, end_key)

(-∞, +∞)
Sorted Map

# PD - Overview

- Meta data management
- Load balance management

# TiKV - Multi-Raft

Multiple raft groups in the cluster, one group for each region.

# TiKV - Horizontal Scale



Node B

Node A

Region 1*

Region 2

Region 3

Add Replica

Node E

Region 1^

Region 2

Region 2

Region 3

Node C

Region 1

Region 3

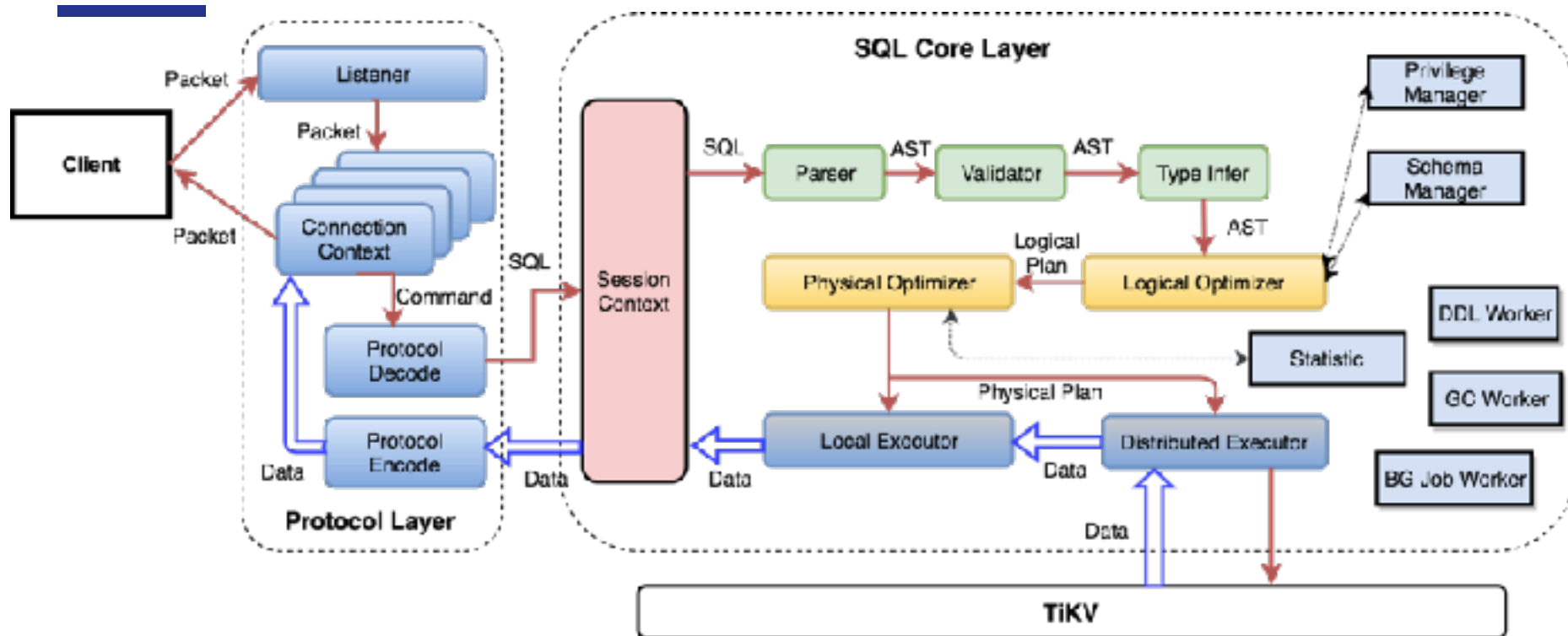Node D

Three steps to move a leader replica
- Transfer Leader
- Add Replica
- Remove Replica

# SQL Layer

# Example - SQL

**CREATE TABLE t (c1 INT, c2 VARCHAR(32), INDEX idx1 (c1));**

**SELECT COUNT(c1) FROM t WHERE c1 > 10 AND c2 = "gopherchina";**

# Example - Logical Plan

# Example - Physical Plan

# Challenges of distributed ACID database?

- Distributed Database is very complex

- Lots of RPC work

- Keep high performance

- Tons of data

- Huge amount of OLTP queries

- Very complex OLAP queries

- External Consistency

- SQL is much more complex than KV

# Why TiDB choose Golang?

- Easy-learning

- Productivity

- Concurrency

- Easy to trace bugs and profile

- Standard libraries and tools

- Tolerant GC latency

- Good performance

- Quick improvement

# Go in TiDB

- More than 160K lines of Go code and 138 contributors.

```
github.com/AlDanial/cloc v 1.72  T=4.00 s (137.4 files/s, 52274.8 lines/s)
-------------------------------------------------------------------------------
Language                   files          blank        comment           code
-------------------------------------------------------------------------------
Go                           523          17655          18761         163173
yacc                           1            389            220           6236
XML                            7              0              0            999
JSON                           1              0              0            502
Markdown                       9            238              0            470
YAML                           3              6              4            227
make                           1             38              2            132
Bourne Shell                   3             12              5             52
Bourne Again Shell             1              9             18             47
Dockerfile                     1              5              0              8
-------------------------------------------------------------------------------
SUM:                         550          18352          19010         171846
-------------------------------------------------------------------------------
```

# Memory && GC

- Query may touch a huge number of data.

- Memory allocation may cost a lot of time.

- Put pressure on GC worker.

- Degrade the performance of SQL.

- OOM sucks!

- runtime.morestack

# Reduce the Number of Allocation

- Get enough memory in one allocation operation

  ```
  a := []int{1, 2, 3, 4, 5}
  b := []int{}
  // a much better way:
  // b := make([]int, 0, len(a))
  for _, i := range a {
      b = append(b, i)
  }
  ```

# Reuse Object

- Share a stack for all queries in one session

- Introduce a cache in goyacc

- Resource pool

# Reduce runtime.morestack

```go
var growStack = false

func growStack() {
    var buf [16 << 10] /* 16 KB */ byte
    if growStack {

        for i := range buf {
            buf[i] = byte(i)
        }
        groupStack = true
    }
}
```

Not enough!

# Goroutine Pool

tidb/util/goroutine_pool/gp.go:

```go
15    package gp
16
17    import (
18        "sync"
19        "sync/atomic"
20        "time"
21    )
22
23    // Pool is a struct to represent goroutine pool.
24    type Pool struct {
25        head        goroutine
26        tail        *goroutine
27        count       int
28        idleTimeout time.Duration
29        sync.Mutex
30    }
31
32    // goroutine is actually a background goroutine, with a channel binded for communication.
33    type goroutine struct {
34        ch     chan func()
35        next   *goroutine
36        status int32
37    }
38
```

# Row Memory Format

| Id (int) | name(varchar) | score(double) | salary(bigint) |
|:---:|:---:|:---:|:---:|
| 1 | "a" | 1.0 | 10000 |
| 2 | "b" | 1.2 | 20000 |
| 3 | "c" | 5.1 | 30000 |
| 4 | "d" | 7.9 | 4000 |

# Row Memory Format

- How to store the row in memory?
- Union? Json? Protobuf?
- Use Datum to store one Column.

```go
type Datum struct {
 k        byte       // datum kind.
 collation uint8      // collation can hold uint8 values.
 decimal   uint16     // decimal can hold uint16 values.
 length    uint32     // length can hold uint32 values.
 i         int64      // i can hold int64 uint64 float64 values.
 b        []byte      // b can hold string or []byte values.
 x        interface{} // x hold all other types.
}
```

# Row Memory Format

- **What is the disadvantage about Datum?**

  - Use unnecessary memory in every column.

  - Must use assert to get complex types:

    ```
    func (d *Datum) GetMysqlDecimal() *MyDecimal {
               return d.x.(*MyDecimal)
    }
    ```

  - Impossible to do vectorizable serial computation

- So how to store multi types data in golang?

# Apache Arrow

- binary data format

- Array lengths

- Null count

- Null bitmaps

- Offsets buffer

- Values Array

- more details plz see the docs

# Row Memory Format

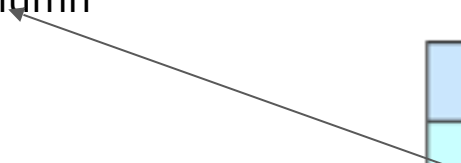- We can eliminate the arrow's offsets array.
- Column store VS Row store.

| id (int) | name(varchar) | score(double) | salary(bigint) |
|----------|---------------|---------------|----------------|
| 1 | "a" | 1.0 | 10000 |
| 2 | "b" | 1.2 | 20000 |
| 3 | "c" | 5.1 | 30000 |
| 4 | "d" | 7.9 | 4000 |

# Chunk

```go
type Chunk struct {
        columns []*column
}

type column struct {
        length      int
        nullCount   int
        nullBitmap  []byte
        offsets     []int32
        data        []byte
        elemBuf     []byte
        ifaces      []interface{}
}
```

| Id (Int) | name(varchar) | score(double) | salary(bigInt) |
|----------|---------------|---------------|----------------|
| 1 | "a" | 1.0 | 10000 |
| 2 | "b" | 1.2 | 20000 |
| 3 | "c" | 5.1 | 30000 |
| 4 | "d" | 7.9 | 4000 |

# Chunk

- **Use unsafe.pointers to get complex types:**

```go
func (r Row) GetMyDecimal(colIdx int) *types.MyDecimal {
    col := r.c.columns[colIdx]
    return (*types.MyDecimal)(unsafe.Pointer(&col.data[r.idx*types.MyDecimalStructSize]))
}

func (r Row) GetUint64(colIdx int) uint64 {
    col := r.c.columns[colIdx]
    return *(*uint64)(unsafe.Pointer(&col.data[r.idx*8]))
}
```

# Chunk

- **Vectorized Execute expressions:**

```go
func VectorizedExecute(ctx context.Context, exprs []Expression, input, output *chunk.Chunk) error {
    sc := ctx.GetSessionVars().StmtCtx
    for colID, expr := range exprs {
        err := evalOneColumn(sc, expr, input, output, colID)
        if err != nil {
            return errors.Trace(err)
        }
    }
    return nil
}

func evalOneColumn(sc *stmtctx.StatementContext, expr Expression, input, output *chunk.Chunk, colID int) (err error) {
    switch fieldType, evalType := expr.GetType(), expr.GetType().EvalType(); evalType {
    case types.ETInt:
        for row := input.Begin(); err == nil && row != input.End(); row = row.Next() {
            err = executeToInt(sc, expr, fieldType, row, output, colID)
        }

    .....
}
```
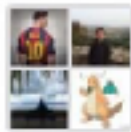
# Thanks!

Contact me:

wink@pingcap.com
www.pingcap.com
Wechat: 82091309

12.16 深圳 Gopher meetup

该二维码7天内(12月22日前)有效，重新进入将更新