



KubeCon



CloudNativeCon

North America 2024

Scale Job Triggering with a Distributed Scheduler

Artur Souza & Cassie Coyle
Diagrid

Table of Contents

- Introduction to Dapr
- Dapr's Architecture
- Actors
- Actor Reminders
- Limitations with Actor Reminders
- Dapr Workflow
- Limitations with Workflow
- Scheduler
- Impact on Actor Reminders and Workflow
- Jobs API
- Conclusion & Future
- Thank you to the contributors that made this possible



 Software Engineer

 @cicoyle

 <https://www.linkedin.com/in/cassie-coyle/>

 Head of Engineering

 @artursouza

 <https://www.linkedin.com/in/barbalho/>



KubeCon



CloudNativeCon

North America 2024

Introduction to Dapr

Developer Challenges

How do I send messages to many services?

How do I ensure messages are sent to particular services?

I just want to trace my calls end-to-end.
What's OpenTelemetry?

How do I implement first-write-wins?

How do I secure access to my data layer?

How do I handle failed calls and perform retries?

I just want to trace my calls end-to-end?

How do I create long-running resilient, stateful services?

How do services discover and call each other securely?



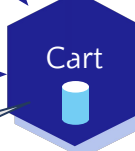
Internet



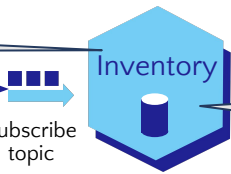
Frontend
API



Email

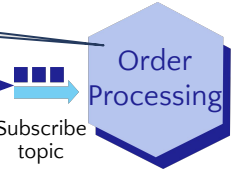


Cart



Inventory

Subscribe
topic



Order
Processing

Subscribe
topic

Publish
event



Checkout

Subscribe
topic

Request-Response
invoke method



Shipping



State
Management



Message
Queue

E-commerce app

Dapr APIs

Application code

Microservices written in



HTTP API

gRPC API



Service-to-
service
invocation



State
management



Publish
and
subscribe



Resource
bindings
and triggers



Actors



Observability



Secrets



Configuration



Distributed
Lock



Workflow



Jobs

Any cloud or edge infrastructure



Virtual or
physical machines



KubeCon



CloudNativeCon

North America 2024



Distributed Application Runtime

Portable, event-driven, runtime for building distributed applications across cloud and edge

dapr.io

The screenshot shows the Dapr website homepage. At the top is a navigation bar with links for Home, Testimonials, Docs, Blog, GitHub, and Discord. On the right of the navigation bar are a star icon with '17,575' and a 'Get Started' button. The main content area features the heading 'APIs for building portable and reliable microservices' and a subheading 'Leverage industry best practices and focus on your application's logic.' Below this is another 'Get Started' button. To the right is a diagram illustrating Dapr's capabilities: a central Dapr icon is connected to three other Dapr icons. The top connection is labeled 'Invoke', the middle 'Store' (with a database icon), and the bottom 'Publish' (with a message icon). The bottom Dapr icon is also connected to a 'Subscribe' label. The footer of the website displays logos for Bosch, Zeiss, Alibaba Cloud, Ignition, Roadwork, 高德地图 (Amap), Legentic, and Man.

Home Testimonials Docs Blog GitHub Discord

☆ Star 17,575 Get Started

APIs for building portable and reliable microservices

Leverage industry best practices and focus on your application's logic.

Get Started

BOSCH ZEISS Alibaba Cloud IGNITION Roadwork 高德地图 LEGENTIC Man

Build connected distributed applications faster

The Distributed Application Runtime (Dapr) provides APIs that simplify microservice connectivity. Whether your communication pattern is service to service invocation or pub/sub messaging, Dapr helps you write resilient and secured microservices.

By letting Dapr's sidecar take care of the complex challenges such as service discovery, message broker integration, encryption, observability, and secret management, you can focus on business logic and keep your code simple.



KubeCon



CloudNativeCon

North America 2024

Dapr's Architecture



KubeCon

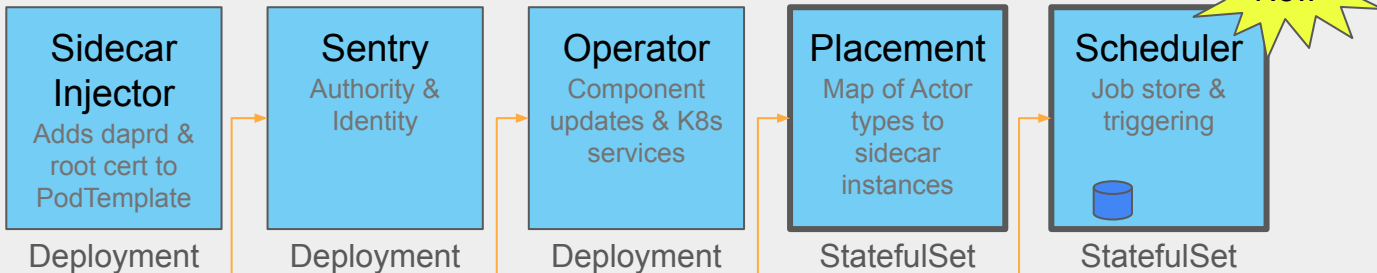


CloudNativeCon

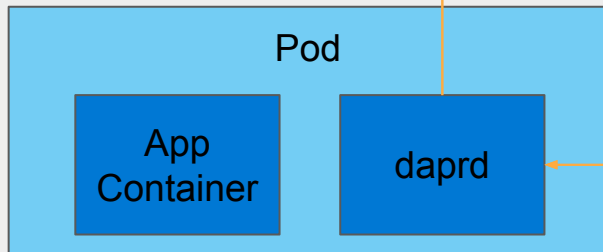
North America 2024

Kubernetes Cluster

dapr-system namespace



App's namespace



Another App's namespace



KubeCon



CloudNativeCon

North America 2024

Actors

POST http://localhost:3500/v1.0/actors/OrderActor/3/method/ship

Machine 1 or Pod 1

Ordering
Service

Actor 1

Actor 2

1



2

Machine 2 or Pod 2

Ordering
Service

Actor 3

Actor 4

3



POST http://localhost/actors/OrderActor/3/method/ship

Video Game
Enemy

X pos
Y pos
Z pos

Difficulty

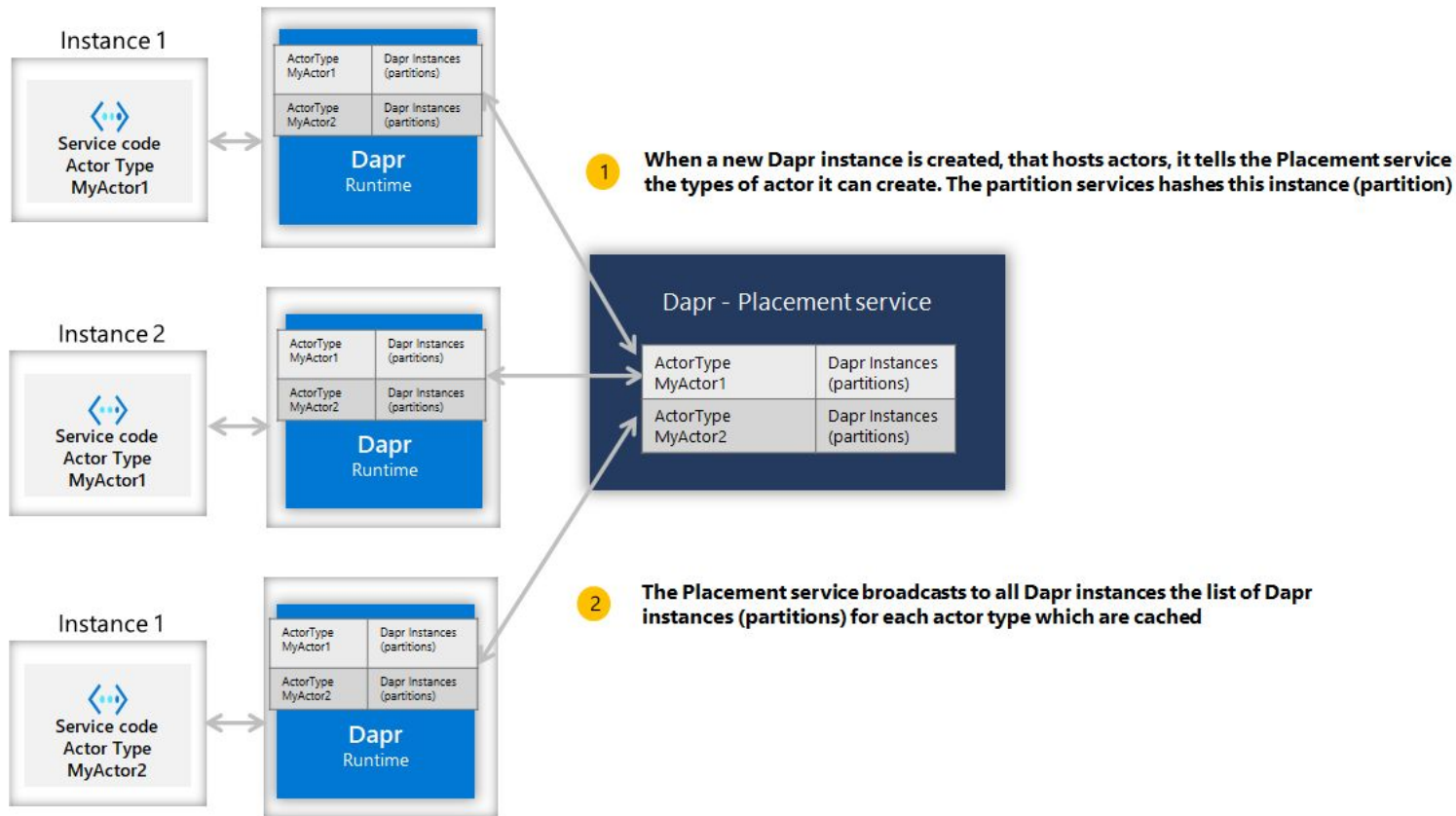
Spawn()

Weapons

Attack()

Host/Pod

Host/Pod





KubeCon



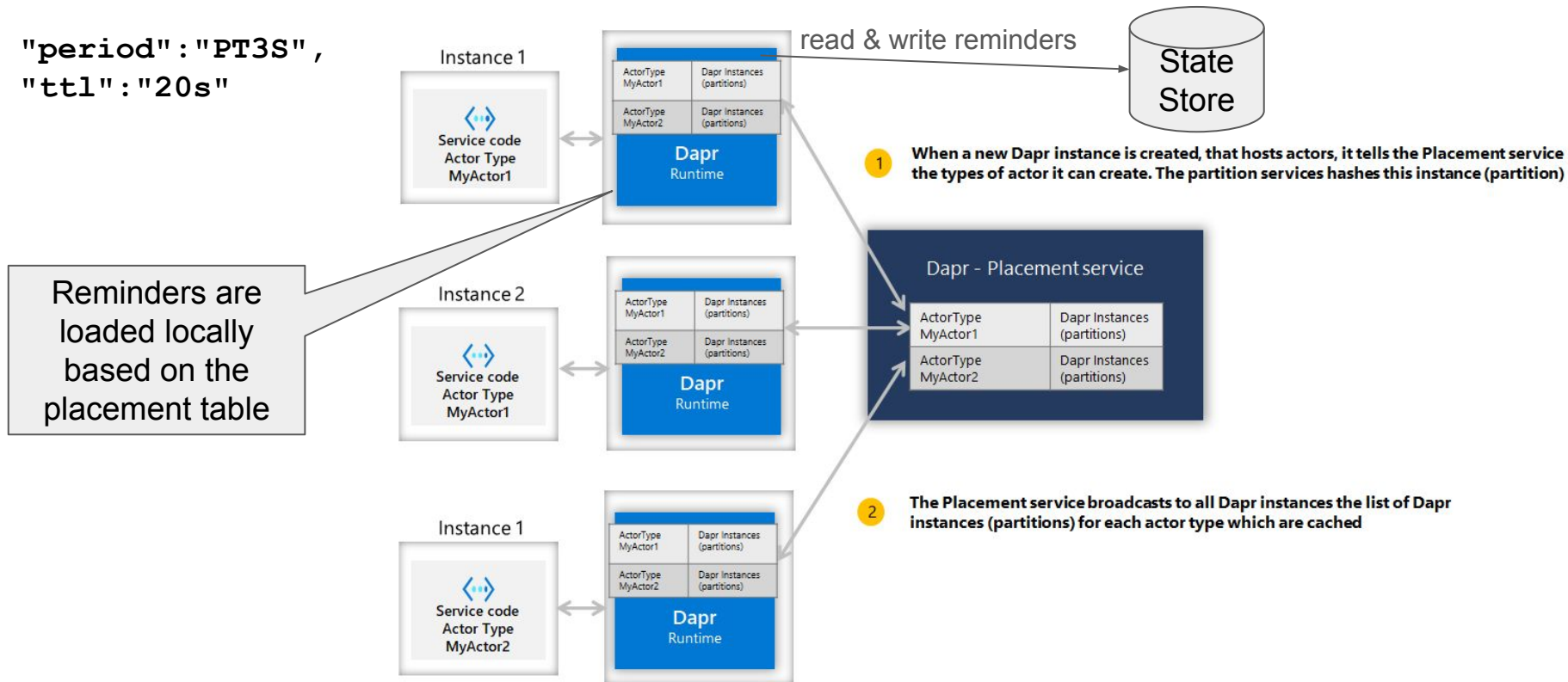
CloudNativeCon

North America 2024

Actor Reminders

POST/PUT `http://localhost:3500/v1.0/actors/<actorType>/<actorId>/reminders/<name>`

```
{
  "period": "PT3S",
  "ttl": "20s"
}
```





Actors||Customer



```
[
  {
    "actorType": "Customer",
    "actorId": "100",
    "name": "paymentReminder",
  },
  {
    "actorType": "Customer",
    "actorId": "200",
    "name": "paymentReminder",
  },
]
```



KubeCon



CloudNativeCon

North America 2024

Key	Value
actors <actor type> metadata	{ "id": <actor metadata identifier>, "actorRemindersMetadata": { "partitionCount": <number of partitions for reminders> } }
actors <actor type> <actor metadata identifier> reminders 1	[<reminder 1-1>, <reminder 1-2>, ... , <reminder 1-n>]
actors <actor type> <actor metadata identifier> reminders 2	[<reminder 1-1>, <reminder 1-2>, ... , <reminder 1-m>]



KubeCon



CloudNativeCon

North America 2024

Key	Value
actors Customer metadata	{ "id": "64d9c7be-8f46-4e1b-9d8a-a95a8aabb43e", "actorRemindersMetadata": { "partitionCount": 2 } }
actors Customer 64d9c7be-8f46-4e1b-9d8a-a95a8aabb43e reminders 1	[{...}, {...}, {...}]
actors Customer 64d9c7be-8f46-4e1b-9d8a-a95a8aabb43e reminders 2	[{...}, {...}, {...}, {...}]



KubeCon



CloudNativeCon

North America 2024

Limitations with Actor Reminders



KubeCon



CloudNativeCon

North America 2024

- Low throughput to register or delete reminders: ~45 tps
- Cannot scale throughput horizontally or vertically
- Limited number of reminders that can be registered: ~1,000 practical limit
- Rebalance required when application pods go up or down



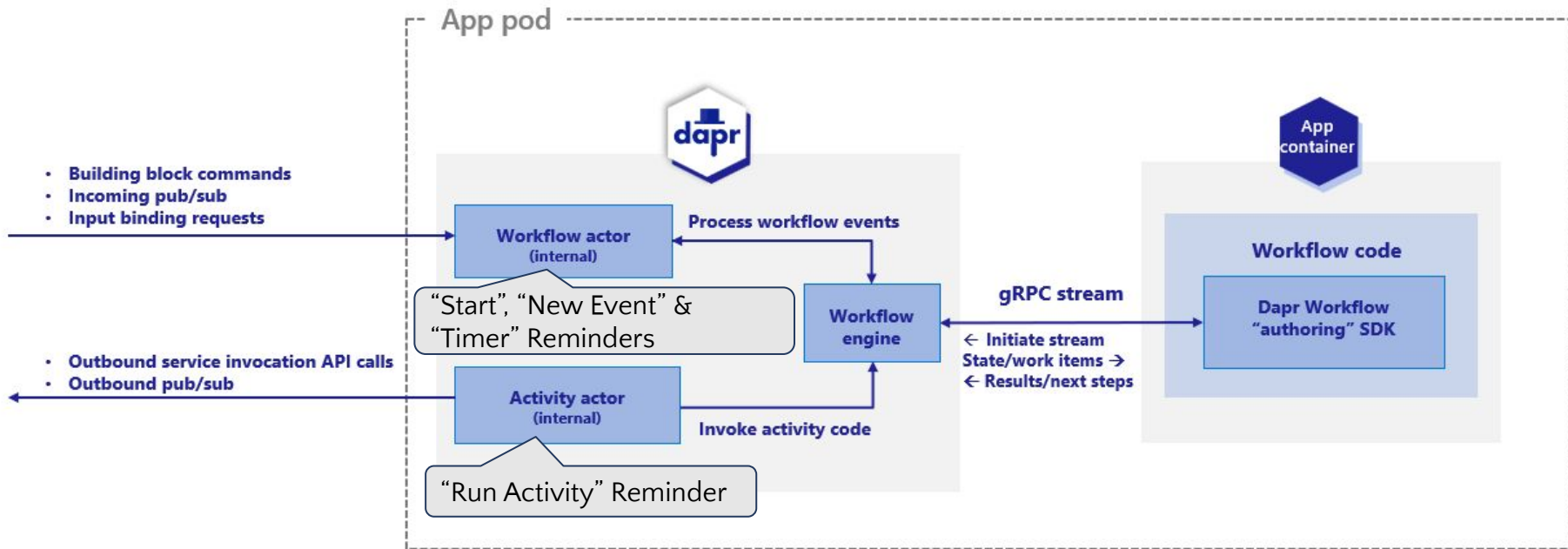
KubeCon



CloudNativeCon

North America 2024

Dapr Workflow





KubeCon



CloudNativeCon

North America 2024



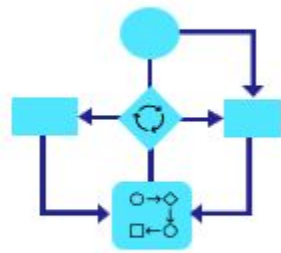
1. Get work item stream



2. Fetch work items



3. Send workflow results





KubeCon



CloudNativeCon

North America 2024

Limitations with Workflows



KubeCon



CloudNativeCon

North America 2024

Practical limit of concurrent workflow activities: ~100

Practical limit of concurrent workflows: 2



KubeCon



CloudNativeCon

North America 2024

Dapr Scheduler Service



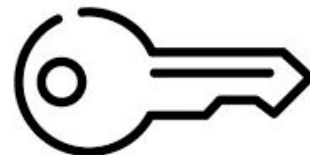
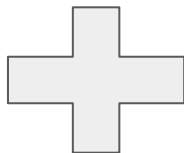
Dapr Scheduler Service

- New control plane service
 - Deployed by default with *dapr init* CLI
 - Run in single instance or in HA mode
- Capabilities
 - Stores jobs to be triggered at some point in the future
 - Guarantees that a job is triggered by one Scheduler
- Implementation
 - Embedded etcd database
 - Internal cron scheduling library



- Job orchestrator, not executor
- At least once job execution
- Bias towards durability and horizontal scaling over clock-time precision
- Generic for multi-purpose job usage

Design Decisions





KubeCon

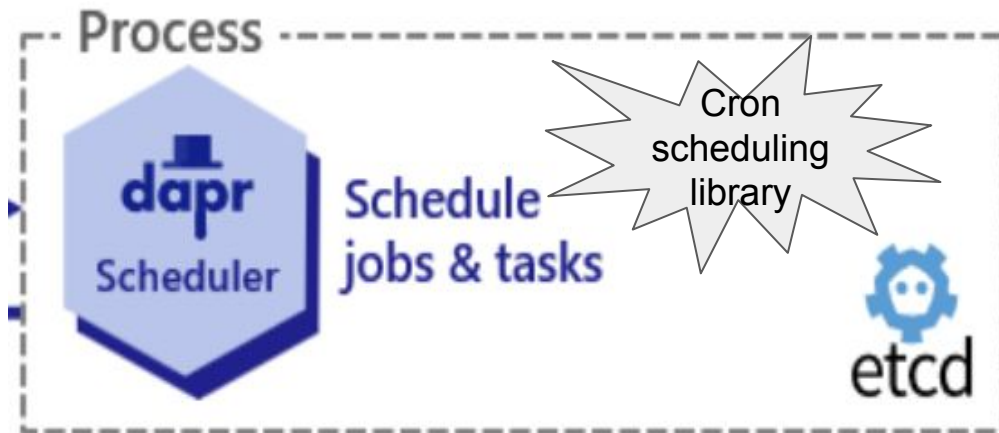


CloudNativeCon

North America 2024

How Does it Work?



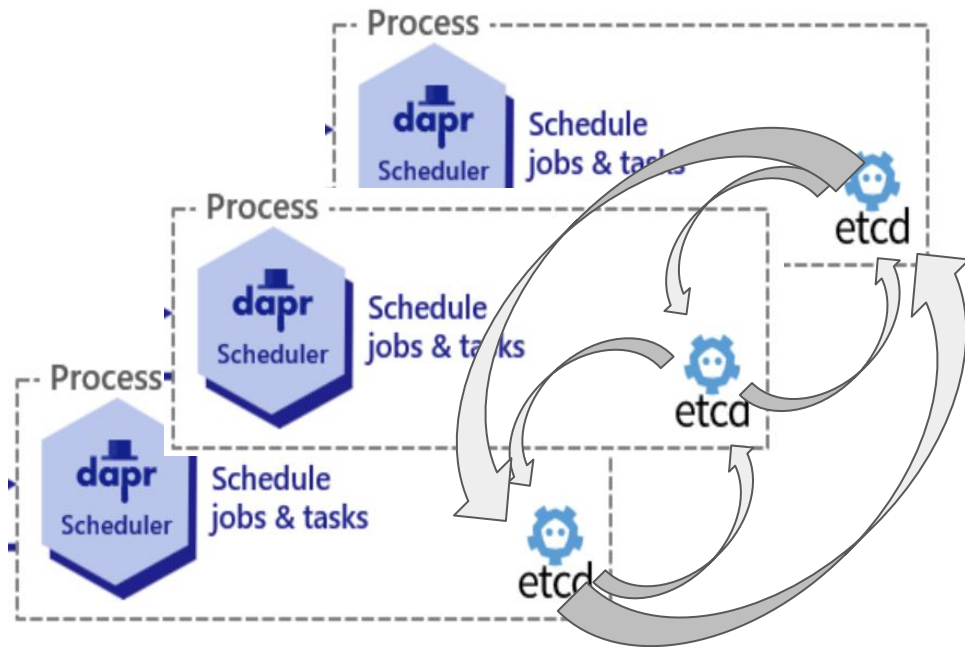


Scheduler Server

```
err = concurrency.NewRunnerManager(  
    ...  
    func(ctx context.Context) error {  
        ...  
        server, serr := server.New(server.Options{  
            Port:             opts.Port,  
            ListenAddress:    opts.ListenAddress,  
            Mode:             modes.DaprMode(opts.Mode),  
            Security:         secHandler,  
            Healthz:          healthz,  
            DataDir:          opts.EtcdDataDir,  
            ReplicaCount:     opts.ReplicaCount,  
            ReplicaID:        opts.ReplicaID,  
            KubeConfig:        opts.KubeConfig,  
            EtcdID:           opts.ID,  
            EtcdInitialPeers:  opts.EtcdInitialPeers,  
            EtcdClientPorts:  opts.EtcdClientPorts,  
            ... (etcd config)  
        })  
        ...  
        return server.Run(ctx)  
    }).Run(ctx)
```


Embedded etcd

- Distributed KV store
- State management of jobs
- Data consistency



Start Embedded etcd

```
import "go.etcd.io/etcd/server/v3/embed"  
  
...  
etcd, err := embed.StartEtcd(c.config)  
if err != nil {  
    return err  
}  
defer etcd.Close()  
  
select {  
case <-etcd.Server.ReadyNotify():  
    log.Info("Etcd server is ready!")  
case <-ctx.Done():  
    return ctx.Err()  
}
```

etcd Data: Leadership

etcd data: Replicated across all instances. The same data for all instances.

dapr/leadership/0	3
dapr/leadership/1	3
dapr/leadership/2	3

etcd Data: Jobs API

etcd data: Replicated across all instances. The same data for all instances.

dapr/jobs/app namespace appid jobid	val

etcd Data: Actor Reminders



KubeCon



CloudNativeCon

North America 2024

etcd data: Replicated across all instances. The same data for all instances.

dapr/jobs/actorreminder default myactortype myactorid remindermethod	val
--	-----

etcd Data: Workflow Actor Reminders

etcd data: Replicated across all instances. The same data for all instances.

<code>dapr/jobs/actorreminder default dapr.internal.default.wf-app.workflow wf-actorid start-eqqANOKQ</code>	val
--	-----

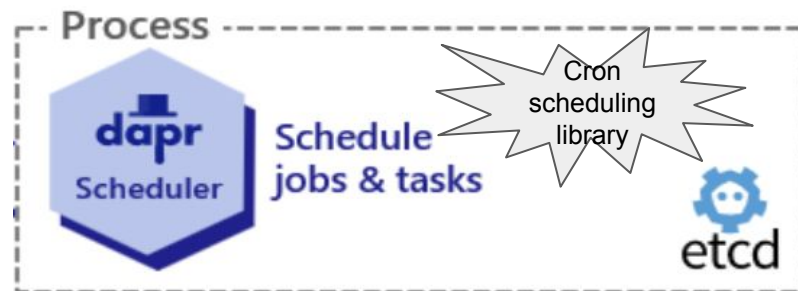
<code>dapr/jobs/actorreminder default dapr.internal.default.wf-app.activity wf-actorid::0::1 run-activity</code>	val
---	-----

<code>dapr/jobs/actorreminder default dapr.internal.default.wf-app.workflow wf-actorid new-event-auiNElfw</code>	val
--	-----

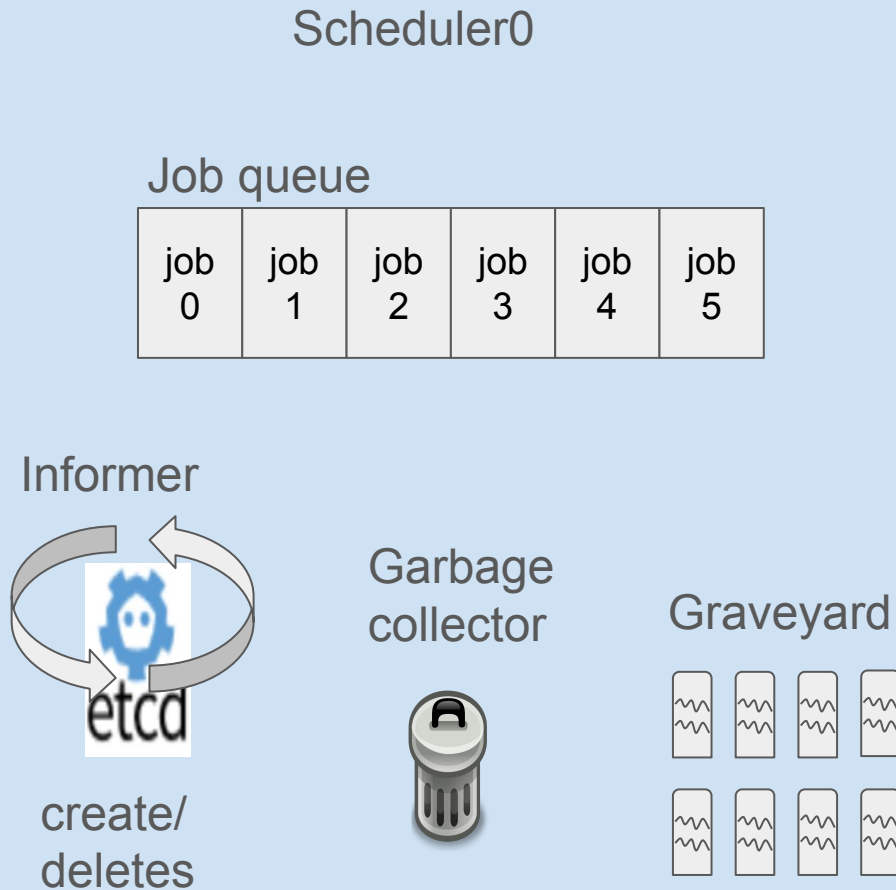
etcd data: Replicated across all instances. The same data for all instances.	
dapr/leadership/0	3
dapr/leadership/1	3
dapr/leadership/2	3
dapr/jobs/app namespace appid jobid	val
dapr/jobs/actorreminder default myactortype myactorid remindermethod	val
dapr/jobs/actorreminder default dapr.interna l.default.wf-app.workflow wf-actorid start-e qqANOKQ	val
dapr/jobs/actorreminder default dapr.interna l.default.wf-app.activity wf-actorid::0::1 r un-activity	val
dapr/jobs/actorreminder default dapr.interna l.default.wf-app.workflow wf-actorid new-eve nt-auiNElfw	val

Internal Cron Scheduling Library

- Enable scalable, distributed job management
- Dynamic job partition leadership coordination
- Enables load distribution



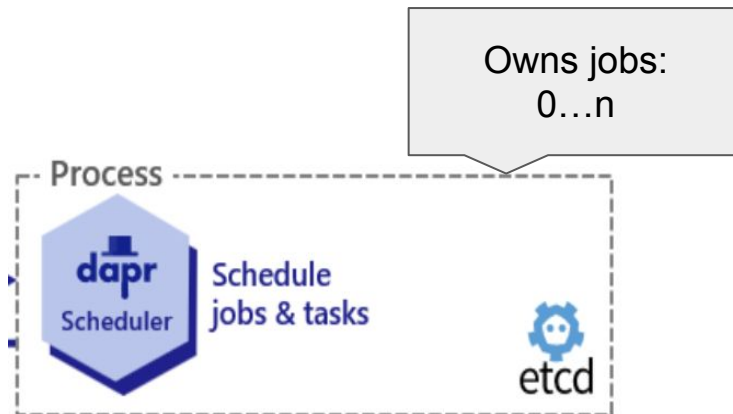
- Queue
 - Manages the scheduling and triggering of jobs
- Counter
 - Track state of triggered jobs over time
- Leadership
 - Job ownership
- Informer
 - Watches for changes in the job keyspace
- Graveyard
 - Track & discard keys
- Garbage Collector
 - Bulk delete keys



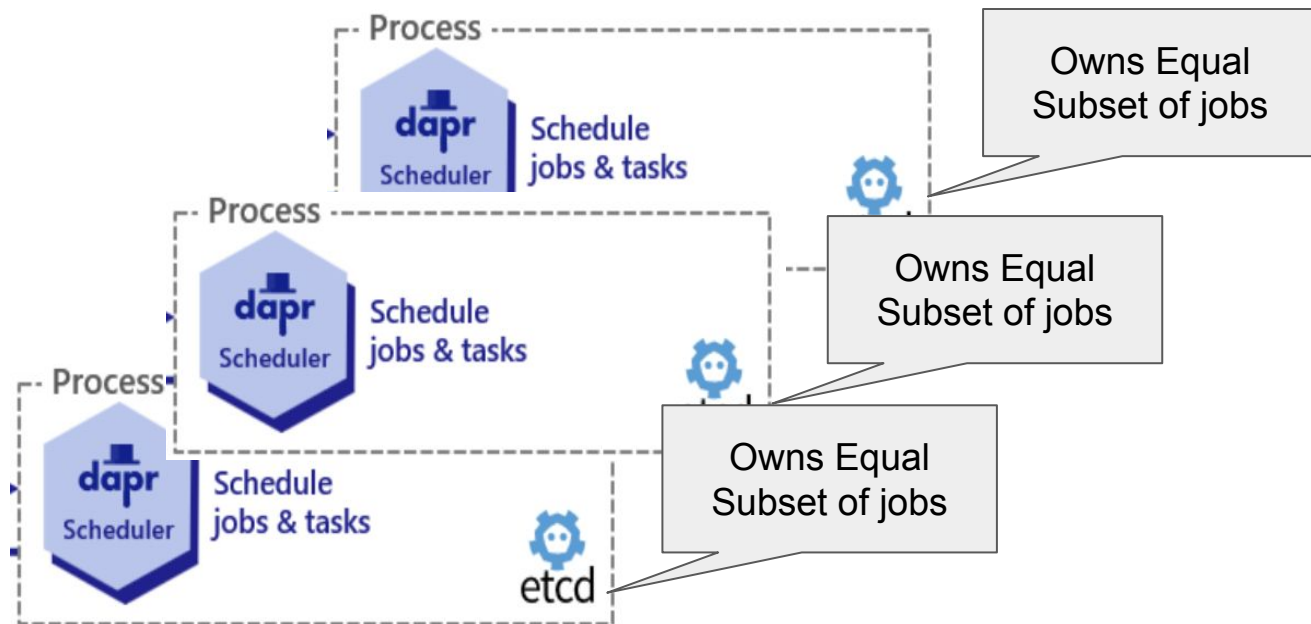
```
import "github.com/diagridio/go-etcd-cron/cron"
...
cron, err := cron.New(cron.Options{
    Client:      client,
    Namespace:   "dapr",
    PartitionID: c.replicaID,
    PartitionTotal: c.replicaCount,
    TriggerFn:   c.triggerJob,
    ReplicaData: replicaData,
})

go cron.Run(context.Background)
```

Cron Library: Ownership model



Cron Library: Ownership model



```
partitionID%m.totalPartitions
```

- Data persistence & replication with etcd
- Dynamic job partition leadership & job distribution
- Jobs are always triggered using the persisted counter for catch-up
- Failure Policy & staging queue



KubeCon

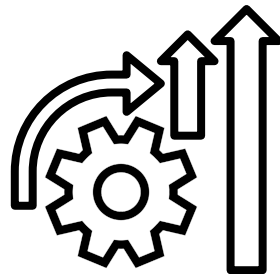


CloudNativeCon

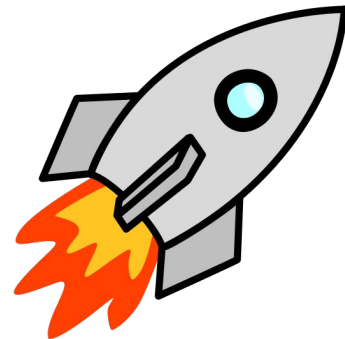
North America 2024

Impact on Actor Reminders and Workflow

- Performance gains in HA mode (3 Scheduler instances)
 - Schedules 50,000 actor reminders with an average trigger QPS of 4,582
 - At least a 10x improvement with stable QPS
 - Direct API invocations achieves up to 35,000 QPS
- Drastic improvements over Dapr v1.13 while creating actor reminders
 - Dapr v1.13 QPS: 50
 - Dapr v1.14 QPS: 4,000 (an 80x increase) with Scheduler



- Parallel Workflow Testing:
 - Max Concurrent Count (60-90):
 - Performance Improvement: 71%
 - Existing Reminder System: Drops by 44%
- High Scale Testing:
 - Max Concurrent Workflows: 350
 - Iterations: 1400
 - Performance Improvement: 50% higher than existing reminder system
- Scale to millions of reminders





KubeCon



CloudNativeCon

North America 2024

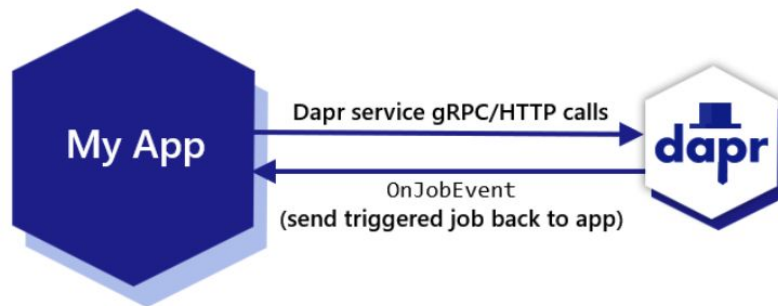
Dapr Jobs API



- Alpha API
- Schedule/Get/Delete

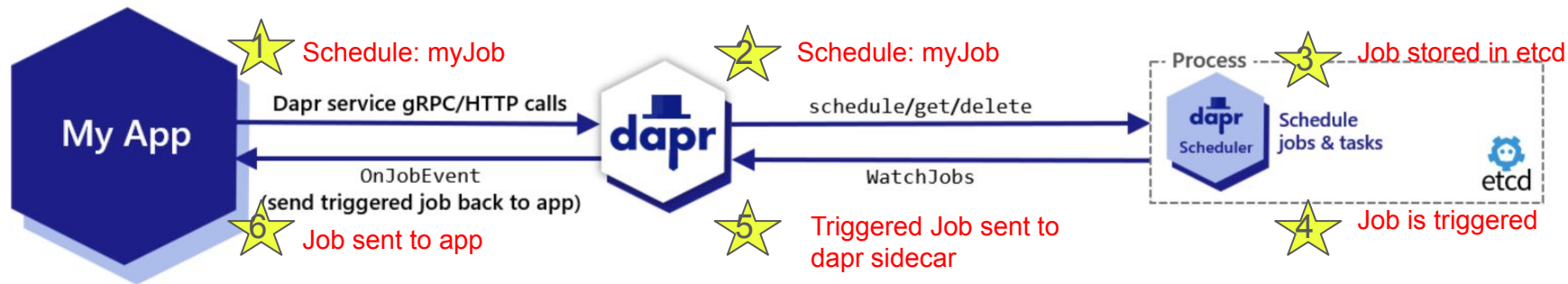
```
curl -X POST \  
  http://localhost:3500/v1.0-alpha1/jobs/test \  
 \  
  -H "Content-Type: application/json" \  
  -d '{  
      "data": "cassie",  
      "dueTime": "3s"  
  }'
```

```
message Job {  
  
    optional string schedule = 1;  
  
    optional uint32 repeats = 2;  
  
    optional string due_time = 3;  
  
    optional string ttl = 4;  
  
    google.protobuf.Any data = 5;  
  
    optional FailurePolicy failure_policy = 6;  
  
}
```



Jobs API + Scheduler Service

- Schedule jobs to be executed at some point in the future



- Delayed PubSub
- Scheduled Service Invocation
- Auto Scale Scheduler
- Optionally Storing Job Data Separately
- CRD job creation

Thank you for your contributions to Dapr Scheduler



KubeCon



CloudNativeCon

North America 2024



Cassie



Elena



Josh



Artur



Yaron