

=====

- شكل الـ App في الـ Legacy Server .. كذا مشكلة منها low utilization of resources بسبب ان منقدرش نحط اكتر من app لهم dependencies مختلفة عن بعض .. وقت عشان اقوم physical server بالتالي low scalability.
- نقلنا للـ vm و containers .. مش بحجز resources في الـ container عشان هو في نظر الـ kernel حيكون process بتاخذ الـ resources اللي محتاجها.
- كمان الـ os بيكون مرة واحدة في الـ container .. لإنه هي process بتيجي لها signal انها تشتغل.
- كمان الـ automation في الـ container وفر لي مبدأ الـ CI/CD اللي بقى لا غنى عنه في عالم الـ microservices .. محتاج fast scaling و الـ fast booting time.
- محتاج الـ Life Cycle Management .. محتاج tool تتأكد إن الـ containers are healthy and running .. محتاج كمان tool تكون مسؤولة عن الـ Containers Scheduling .. لو في ماشين خلاص بقت full utilized .. فلو في container محتاج يـ scheduled فيروح لماشين less utilized .. برضه تكون مسؤولة عن الـ Horizontal / Vertical Scaling .. كمان Health check Monitor للـ Apps ... مش مجرد Start / Stop للـ Containers .. لو الـ Pod محتاجة gpu node .. فالـ tool هي المسؤولة عن الـ required node assignment.
- كمان الـ k8s مش بيـ manage pods فقط ... كمان بيـ manage nodes .. لو عندي node عاوز اخرجها من الـ cluster لـ failure أو maintenance مثلاً ... فـ k8s مسؤول عن الـ running pods on current node to the healthy available nodes ... عن طريق الـ Kublete اللي بيـ report الـ node health/status باستمرار للـ API Server.
- بعد ما الـ node بتراجع بتكون monitored by k8s وبيتعمل عليها some checks ولما تـ passed .. بتبدأ في استقبال / scheduling الـ Pods.

:K8S Features

=====

- 01-Services Discovery: بتقوم مقام الـ stable frontend لكل الـ pods اللي نقدر نقول عليها انها - not stable backend - على مستوى نفس الـ application.
- 02-Load Balancing: لو عاوز اعمل expose الـ service بره الـ cluster عشان تكون managed by a dedicated load balancer. بتعريف الـ ClusterIP والـ NodePort. ودول كده نوعين من الـ Endpoints اللي تخلينا قادرين نتعامل مع الـ pods <<< لكن برضه k8s بيقدر يـ balance loads بشكل efficient وسريع جداً لأن كل الـ traffic جوه الـ cluster .. بتشوف الـ pod الجاهز لاستقبال الـ traffic وتبدأ تبعته ليها.

<<<< Service Object من الـ >>>>

- 03-Self-Healing: عن طريق طريقتين بيقدر يكون متوفر لنا الـ Reliability .. عن طريق طريقتين الـ Readiness probe and Lifeness Probe. الـ Readiness مش بس بعد ما الـ pod بتقوم .. لكن كمان على طول الـ Life Cycle << طول ما الـ pod is ready فيتقدر تستقبل الـ traffic.
- 04-Application Hosting: الـ Developer بيكون Abstracted عن أي حاجة Infrastructure وبيكون تركيزه فقط على الـ Application.

:Cluster Archeticture

=====

- ايه مكوناته؟ .. k8s عبارة عن distributed platform بتتكون من Master as Server و Worker as Client ...
- الماستر هو اللي بيعت ال instructions لل Worker nodes <<< بيتكون من API Server وده ال Central Communication / Gateway اللي بتهدل ال Requests بين باقي مكونات ال master node أو حتى ال worker node من نفس ال cluster أو nodes from other cluster ،،، يجي بعد كده ال Controll Manager المسؤول عن وصول ال Current Cluster State لل Desired Cluster State .. بيفضل يقارن بين ال Current و ال Desired ويتأكد انهم Matched. و كمان دايماً باستمرار بيكلم ال API Server عشان لو في Task مطلوب تنفيذه يقوم بده.
- وعندنا كمان ال Scheduler: هو المسؤول عن pod assignment for suitable node. عن طريق factors مختلفة عشان يختار أفضل node يحط ال pod عليها.

بالنسبة لل Worker Node:

- عندي أول حاجة ال Kubelet ده نقدر نعتبره ال agent اللي بيكون موجود وشغال في كل Worker Node وهو اللي بيسجلها عند ال Master Node API Server وبيبلغ باستمرار ال Worker Node Status ليه .. و كمان بتستنى منه أي instruction اللي المفروض تنفذها.

:Create Object Process

=====

- 1- أو حاجة بيعت لل Master Node API Server ال Request اللي بيقول له فيه انه محتاج يعمل Desired State لـ Deployment ولتكن Nginx مثلاً.
- 2- ال API Server بعد ما يستلم الطلب بيبدأ يعمل بعض ال Checks اللي من خلالها بي Authenticate ال Client الأول "logging" << وبعد كده اشوف اذا كان ال Request مسموح لليوزر انه يعمل ولا لأ .. إلخ.
- 3- يجي دور ال ETCD بعد ما ال Checks تكون Passed .. ال API Server بيخزن ال Deployment Object في ال ETCD Database
- 4- ال ETCD بدوره حيلغ ال API Server انه خزن ال Object بنجاح.
- 5- ال API Server بيلغ ال User إن ال Object Created.
- 6- ال Controller Manager زي باقي ال Components بي listen على ال API Server وحيلاقي تاسك / لوجيك إن مطلوب يكون في Deployment مطلوب انها تكون موجودة في ال cluster وهي مش موجودة .. فال CM حيلاقي انه ال Desired State != Current State <<< فاللي حيحصل ان ال CM حيعمل Create ال Deployment وبعده ال ReplicaSet Controller حيد Create Pods بعدد حسب المطلوب منه في ال Deployment Object << كل ال Objects دي بترجع باستمرار على ال API Server و ال API Server بيخزنها باستمرار على ال ETCD. ال Pods اتعملت ولكن لسه not assigned على أي node <<< لسه ال Status بتاعتها Pending.
- 7- ال Scheduler زيه زي ال CM وهو بيسمع ال API Server .. فحيلاقي مبعوت له Nodes بـ Pending Status <<< فحينفذ اللوجيك بتاعه ويشوف ال nodes ويبدأ ي assign ال pod object requirements ويختار له ال node اللي ت match them ... هنا ال Sched مش مطلوب منه تشغيل ال pod .. كل اللي مطلوب منه ي update ال node specs لكل pod في ال ETCD عن طريق ال API Server برضه.
- 8- ال Kubelete باستمرار برضه سامعة ال API Server .. فحتلاقي لوجيك مطلوب منها بان في 3 pods مثلاً مطلوب 2 منهم يكونوا على node 1 و 1 منهم على node 2 فهنا ال Kubelete حتشغل ال pods ولكن مش هي اللي حتعمل ال Container.

- 9- لكن الـ Container Runtime هي اللي تحتسلم Request من الـ Kubelete وهي اللي حتعمل الـ Container وبكده الـ Deployment اتعملت والـ Pods اشتغلت.

:Master Components

=====

- الـ communication بين الـ API Server وبين أي Component بيتم عن طريق Certificate / متشفر زي الـ PKI.
- من ضمن الـ pod specs اللي بيحدد على اساسها الـ scheduler حيت assign كل pod على انه node .. ان الـ pod محتاج affinity أو anit-affinity يعني هل مجموعة الـ pods شرط انها تشتغل مع بعض على نفس الـ node ولا العكس على nodes مختلفة.
- Etcd مش مقبول انها تكون down تحت أي ظرف لإن ده معناه ان الـ Cluster كله down. بالتالي من اهميتها انا محتاج اضمن ان يكون معايا backup باستمرار منها .. عشان اقدر restore state of cluster .. في كل master node بيكون معايا جواها etcd instance وده موديل بيكون اسمه stacked etcd و الـ api server بيكلم الـ instance ده. أو ممكن نجمع الـ etcd instanced في external cluster والـ api server بيتكلم معاها. ده افضل عشان اضمن High Availability.
- كام master node بيكون على حسب حجم الـ traffic من cluster traffic اللي بييجي على api server لو الـ user محتاج يعمل object أو يـ scale deployment يعمل service ...
- لكن الـ app traffic بيروح للـ ingress controller ومنه بيوجه الـ traffic ناحية الـ worker nodes وبالتالي ده محسوب معانا كـ cluster traffic .. لكن لو في update / scale request محتاج يعمل الـ app فساعتها ده مش حيتم لو الـ master node واقعين كلهم.
- كام etcd ؟ ... يفضل يكون الحد الأدنى 3 والحد الأقصى مش مفتوح .. لسبب إن الـ writing بيكون sequential بسبب ان مسموح لـ etcd node leader فقط والـ follower بيد sync data كل واحد مع اللي قبلها بشكل sequential as chain فكده العدد الكبير منهم حيسبب latency عقبال ما الـ update يحصل على كل الـ etcd nodes.
- بالنسبة للـ Request Validation اللي بييجي من الـ user: هل الـ user له access / authenticated على الـ cluster .. لو ليه فحيشوف لو هو authorized / has permission انه يعمل له الـ request ... بعد كده بنتيجي مرحلة Mutating Admission Controller ... هنا بيتغير شكل الـ object تحسباً لو في bug في الـ object specs أو حاجة الـ user مش محدد تقدر تخلي الـ pod بشكلها الحالي تاخذ كل الـ resources اللي في الـ cluster .. في المثال ده المسؤول عن تفادي ان ده يحصل حاجة اسمها limit range admission controller.
- بعد كده تيجي مرحلة الـ Schema Validation ... بيكون للـ object بعد ما اتعدل من الـ M A C قبل ما يسجله في الـ etcd database.
- بنتيجي بعد كده اخر مرحلة Validation Admission Controller ... هنا ده final validation اني بحط policy زي ان اليوزر مش مسموح له يـ deploy public image .. لازم تكون من الـ internal registry ... أو انه ممنوع من استخدام الـ latest tag ... أو الـ deployment تكون باستخدام deployment pipeline زي الـ Jenkins ... لازم تكون باستخدام authorized pipeline.
- لما كل ده يتم ... تقدر نخط الـ object في الـ etcd database.

:Worker Components

=====

- Kubelet تعتبر الـ agent التي بتكون موجودة في كل node ودي بتكون مسؤولة عن تسجيل / register و تـ report the status الـ node في الـ API Server من readiness & liveness & utilization.
- بيكون في interface بين الـ kubelet وبين الـ "CRI" Container Runtime بتبعت فيه تعليمات وشكل الـ container المطلوب.
- Kube Proxy موجود في كل worker node ودوره هو service management .. عن طريق انه برضه بيشوف الـ API Server باستمرار .. وأول ما يلاقي ان في service اتعملت .. فهو اللي بيعمل لها rule بـ virtual ip خاص بيها في الـ ip tables ... نفس الكلام بالنسبة للـ pods لما بتتعمل لكنه مش بيهندلها هي بشكل مباشر .. بيكون الـ Endpoint controller manager هو اللي أول ما يلاقي pod/s اتعملت بيعمل لها حاجة اسمها Endpoint ... الفكرة اني ضامن ان ده ip قادر انه يستقبل traffic بالتالي يقدر يتعامل معاه الـ kube proxy ويعمل لها rule بالـ Service و الـ Endpoints.