

AUTOSCALING - OVERVIEW

- THERE ARE TWO MAIN WAYS THAT AN APPLICATION CAN SCALE:
 - **VERTICAL SCALING**, ALSO CALLED SCALING UP AND DOWN, MEANS CHANGING THE CAPACITY OF A RESOURCE. FOR EXAMPLE, YOU COULD MOVE AN APPLICATION TO A LARGER VM SIZE. VERTICAL SCALING OFTEN REQUIRES MAKING THE SYSTEM TEMPORARILY UNAVAILABLE WHILE IT IS BEING REDEPLOYED. THEREFORE, IT'S LESS COMMON TO AUTOMATE VERTICAL SCALING.
 - **HORIZONTAL SCALING**, ALSO CALLED SCALING OUT AND IN, MEANS ADDING OR REMOVING INSTANCES OF A RESOURCE. THE APPLICATION CONTINUES RUNNING WITHOUT INTERRUPTION AS NEW RESOURCES ARE PROVISIONED. WHEN THE PROVISIONING PROCESS IS COMPLETE, THE SOLUTION IS DEPLOYED ON THESE ADDITIONAL RESOURCES. IF DEMAND DROPS, THE ADDITIONAL RESOURCES CAN BE SHUT DOWN CLEANLY AND DEALLOCATED.
- ALL 3 CLOUD PLATFORMS SUPPORT HORIZONTAL SCALING

AUTO SCALING - AWS

- AT THE TIME OF INITIAL LAUNCH, AWS DID NOT OFFER AUTOSCALING, BUT THE ABILITY TO PROGRAMMATICALLY CREATE AND TERMINATE INSTANCES GAVE DEVELOPERS THE FLEXIBILITY TO WRITE THEIR OWN CODE FOR AUTOSCALING
- THIRD-PARTY AUTOSCALING SOFTWARE FOR AWS BEGAN APPEARING AROUND APRIL 2008. THESE INCLUDED TOOLS BY SCALR AND RIGHTS SCALE
- ON MAY 18, 2009, AMAZON LAUNCHED ITS OWN AUTOSCALING FEATURE ALONG WITH ELASTIC LOAD BALANCING
- ON-DEMAND VIDEO PROVIDER NETFLIX DOCUMENTED THEIR USE OF AUTOSCALING WITH AMAZON WEB SERVICES TO MEET THEIR HIGHLY VARIABLE CONSUMER NEEDS

AUTOSCALING - AWS

- VARIOUS BEST PRACTICE GUIDES FOR AWS USE SUGGEST USING ITS AUTOSCALING FEATURE EVEN IN CASES WHERE THE LOAD IS NOT VARIABLE
- THAT IS BECAUSE AUTOSCALING OFFERS TWO OTHER ADVANTAGES:
 - AUTOMATIC REPLACEMENT OF ANY INSTANCES THAT BECOME UNHEALTHY FOR ANY REASON (SUCH AS HARDWARE FAILURE, NETWORK FAILURE, OR APPLICATION ERROR)
 - AUTOMATIC REPLACEMENT OF SPOT INSTANCES THAT GET INTERRUPTED FOR PRICE OR CAPACITY REASONS, MAKING IT MORE FEASIBLE TO USE SPOT INSTANCES FOR PRODUCTION PURPOSES

AUTOSCALING - AZURE

- ON JUNE 27, 2013, MICROSOFT ANNOUNCED THAT IT WAS ADDING AUTOSCALING SUPPORT TO ITS WINDOWS AZURE CLOUD COMPUTING PLATFORM
- CONDITIONS CAN BE SET FOR A CLOUD SERVICE WORKER ROLE THAT TRIGGER A SCALE IN OR OUT OPERATION
- THE CONDITIONS FOR THE ROLE CAN BE BASED ON THE CPU, DISK, OR NETWORK LOAD OF THE ROLE
- YOU CAN SCALE AN APPLICATION ONLY WITHIN THE LIMIT OF CORES FOR YOUR SUBSCRIPTION
- TO ENABLE HIGH AVAILABILITY OF YOUR APPLICATION, YOU SHOULD ENSURE THAT IT IS DEPLOYED WITH TWO OR MORE ROLE INSTANCES
- AUTO SCALE ONLY HAPPENS WHEN ALL THE ROLES ARE IN **READY** STATE

AUTOSCALING - AZURE

- YOU CAN CONFIGURE SCALE SETTINGS FOR A ROLE WITH EITHER TWO MODES **MANUAL** OR **AUTOMATIC**
- MANUAL IS AS YOU WOULD EXPECT, YOU SET THE ABSOLUTE COUNT OF INSTANCES
- AUTOMATIC HOWEVER ALLOWS YOU TO SET RULES THAT GOVERN HOW AND BY HOW MUCH YOU SHOULD SCALE
- YOU CAN SCALE
 - ALWAYS
 - RECURRENCE (SET OF DAYS)
 - FIXED DATE

AUTOSCALING - GCP

- ON NOVEMBER 17, 2014, THE GOOGLE COMPUTE ENGINE ANNOUNCED A PUBLIC BETA OF ITS AUTOSCALING FEATURE FOR USE IN GOOGLE CLOUD PLATFORM APPLICATIONS
- AUTOSCALING IS A FEATURE OF MANAGED INSTANCE GROUPS. A MANAGED INSTANCE GROUP IS A POOL OF HOMOGENEOUS INSTANCES, CREATED FROM A COMMON INSTANCE TEMPLATE
- SCALING USING THE FOLLOWING POLICIES:
 - AVERAGE CPU UTILIZATION
 - HTTP LOAD BALANCING SERVING CAPACITY, WHICH CAN BE BASED ON EITHER UTILIZATION OR REQUESTS PER SECOND
 - STACKDRIVER MONITORING METRICS

SERVERLESS COMPUTING - OVERVIEW

- **SERVERLESS COMPUTING** IS A CLOUD **COMPUTING** EXECUTION MODEL IN WHICH THE CLOUD PROVIDER DYNAMICALLY MANAGES THE ALLOCATION OF MACHINE RESOURCES.
- AWS OFFERS THE FOLLOWING PRODUCTS WITHIN IT'S SERVERLESS PLATFORM; AWS LAMBDA, AMAZON API GATEWAY, AMAZON SIMPLE STORAGE SERVICE, AMAZON DYNAMODB, AMAZON SNS, AMAZON SQS, AWS STEP FUNCTIONS, AMAZON KINESIS.
- AZURE OFFERS THE FOLLOWING PRODUCTS WITHIN IT'S SERVERLESS PLATFORM; AZURE FUNCTIONS, AZURE STORAGE, AZURE COSMOS DB, AZURE ACTIVE DIRECTORY, EVENT GRID, SERVICE BUS, LOGIC APPS, AZURE STREAM ANALYTICS, AZURE BOT SERVICE, COGNITIVE SERVICES
- GCP OFFERS THE FOLLOWING PRODUCTS WITHIN IT'S SERVERLESS PLATFORM; APP ENGINE, CLOUD FUNCTIONS, CLOUD DATASTORE, CLOUD STORAGE, CLOUD PUB/SUB, APIGEE, ENDPOINTS, CLOUD DATAFLOW, BIGQUERY, CLOUD ML ENGINE

AWS LAMBDA FEATURES

- SUPPORTED PROGRAMMING LANGUAGES - JAVA, NODE.JS, C#, AND PYTHON CODE, WITH SUPPORT FOR OTHER LANGUAGES COMING IN THE FUTURE.
- AUTOMATIC SCALING - AWS LAMBDA AUTOMATICALLY SCALES TO SUPPORT THE RATE OF INCOMING REQUESTS WITHOUT REQUIRING YOU TO CONFIGURE ANYTHING.
- INTEGRATION WITH OTHER SERVICES - YOU CAN CONFIGURE OTHER RESOURCES SUCH AS S3 BUCKETS, DYNAMODB TABLE, OR KINESIS STREAM TO TRIGGER LAMBDA FUNCTIONS.
- FAULT TOLERANCE - AWS LAMBDA MAINTAINS COMPUTE CAPACITY ACROSS MULTIPLE AVAILABILITY ZONES IN EACH REGION TO HELP PROTECT YOUR CODE AGAINST INDIVIDUAL MACHINE OR DATA CENTER FACILITY FAILURES.
- PRICING MODEL - BILLING IS METERED IN INCREMENTS OF 100 MILLISECONDS, MAKING IT COST-EFFECTIVE AND EASY TO SCALE AUTOMATICALLY FROM A FEW REQUESTS PER DAY TO THOUSANDS PER SECOND.
- CONTINUOUS INTEGRATION - TO AUTOMATE THE DEPLOYMENT PROCESS, YOU CAN USE THE FOLLOWING SERVICES, CODEPIPELINE, CODEBUILD, CODEDEPLOY, AWS CLOUDFORMATION.

AZURE FUNCTIONS

- SUPPORTED PROGRAMMING LANGUAGES - C#, F#, OR JAVASCRIPT.
- AUTOMATIC SCALING - AZURE FUNCTIONS USES A COMPONENT CALLED THE SCALE CONTROLLER TO MONITOR THE RATE OF EVENTS AND DETERMINE WHETHER TO SCALE OUT OR SCALE IN.
- INTEGRATION WITH OTHER SERVICES - THE FOLLOWING SERVICE INTEGRATIONS ARE SUPPORTED. AZURE COSMOS DB, AZURE EVENT HUBS, AZURE EVENT GRID, AZURE MOBILE APPS, AZURE NOTIFICATION HUBS, AZURE SERVICE BUS (QUEUES AND TOPICS), AZURE STORAGE (BLOB, QUEUES, AND TABLES), GITHUB (WEBHOOKS), ON-PREMISES (USING SERVICE BUS), TWILIO (SMS MESSAGES).
- PRICING MODEL - **CONSUMPTION PLAN** : YOU ONLY PAY FOR THE TIME THAT YOUR CODE RUNS, **APP SERVICE PLAN**: IF YOU ARE ALREADY USING APP SERVICE FOR YOUR OTHER APPLICATIONS, YOU CAN USE THE SAME PLAN.
- CONTINUOUS INTEGRATION - CODE YOUR FUNCTIONS RIGHT IN THE PORTAL OR SET UP CONTINUOUS INTEGRATION AND DEPLOY YOUR CODE THROUGH GITHUB, VISUAL STUDIO TEAM SERVICES, AND OTHER SERVICES.

GCP – CLOUD FUNCTIONS

- SUPPORTED PROGRAMMING LANGUAGES – GO, JAVA, .NET, NODE.JS, PHP, PYTHON, RUBY.
- AUTOMATIC SCALING - CLOUD FUNCTIONS AUTOMATICALLY MANAGES AND SCALES UNDERLYING INFRASTRUCTURE WITH THE SIZE OF WORKLOAD.
- INTEGRATION WITH OTHER SERVICES - CLOUD FUNCTIONS ALLOWS YOU TO TRIGGER YOUR CODE FROM GOOGLE CLOUD PLATFORM, FIREBASE, AND GOOGLE ASSISTANT, OR CALL IT DIRECTLY FROM ANY WEB, MOBILE, OR BACKEND APPLICATION VIA HTTP.
- PRICING MODEL - PAY ONLY WHILE YOUR FUNCTION IS EXECUTING, METERED TO THE NEAREST 100 MILLISECONDS, AND PAY NOTHING AFTER YOUR FUNCTION FINISHES.
- CONTINUOUS INTEGRATION - YOU CAN CONFIGURE A CONTINUOUS INTEGRATION AND DEPLOYMENT (CI/CD) PLATFORM SUCH AS CLOUD CONTAINER BUILDER TO RUN YOUR EXISTING CLOUD FUNCTIONS TESTS ON AN ONGOING BASIS.
- LOCAL EMULATOR - THE CLOUD FUNCTIONS EMULATOR ALLOWS YOU TO **DEPLOY**, **RUN**, AND **DEBUG** YOUR CLOUD FUNCTIONS ON YOUR LOCAL MACHINE BEFORE DEPLOYING THEM TO THE PRODUCTION CLOUD FUNCTIONS SERVICE.

WHAT IS LOAD BALANCING

- A LOAD BALANCER IS A DEVICE THAT ACTS AS A REVERSE PROXY AND DISTRIBUTES NETWORK OR APPLICATION TRAFFIC ACROSS A NUMBER OF SERVERS. LOAD BALANCERS ARE USED TO INCREASE CAPACITY (CONCURRENT USERS) AND RELIABILITY OF APPLICATIONS.
- LOAD BALANCING AIMS TO OPTIMIZE RESOURCE USE, MAXIMIZE THROUGHPUT, MINIMIZE RESPONSE TIME, AND AVOID OVERLOAD OF ANY SINGLE RESOURCE
- USING MULTIPLE COMPONENTS WITH LOAD BALANCING INSTEAD OF A SINGLE COMPONENT MAY INCREASE RELIABILITY AND AVAILABILITY THROUGH REDUNDANCY.
- LOAD BALANCING USUALLY INVOLVES DEDICATED SOFTWARE OR HARDWARE, SUCH AS A MULTILAYER SWITCH OR A DOMAIN NAME SYSTEM SERVER PROCESS.

AWS - ELASTIC LOAD BALANCING

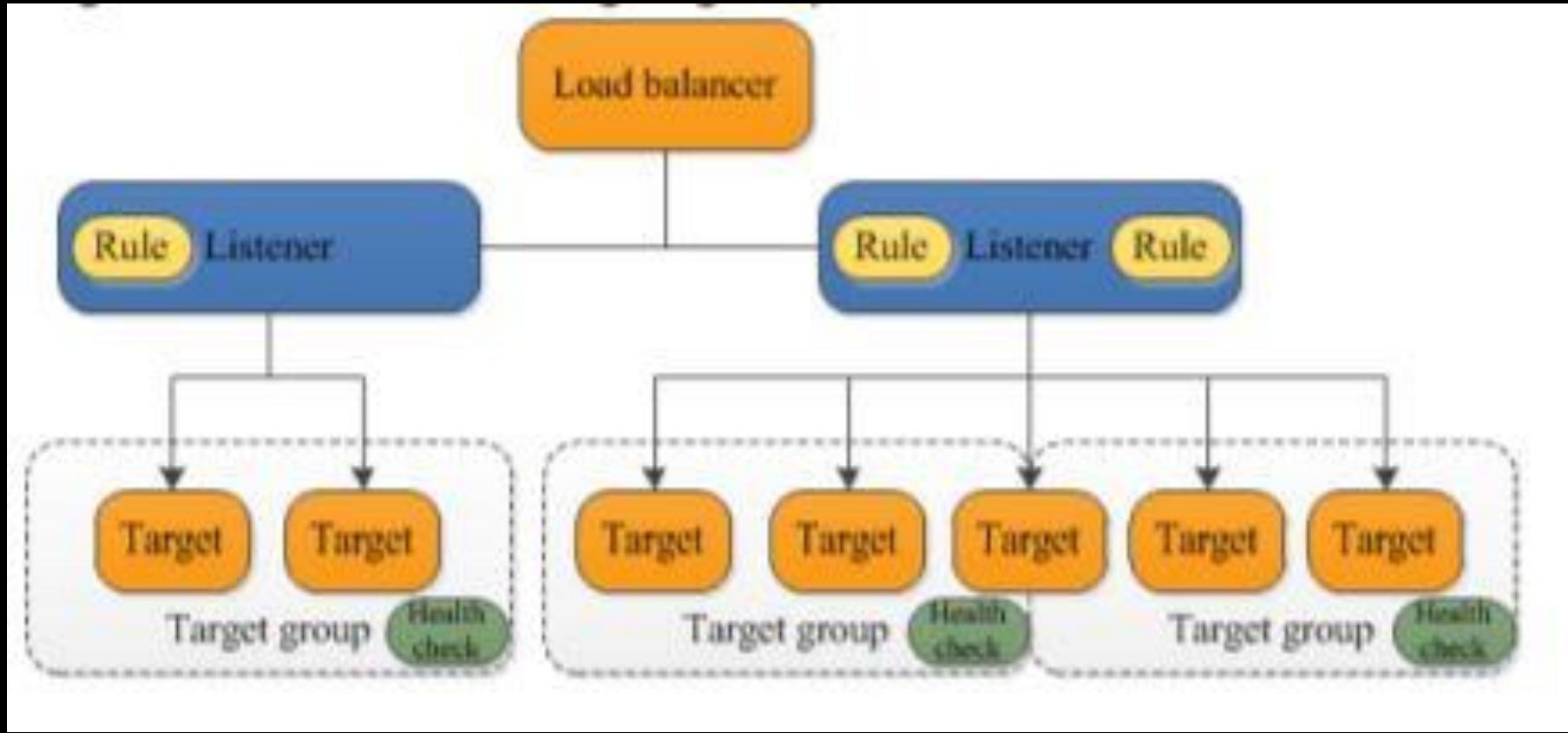
- ❖ ELASTIC LOAD BALANCING DISTRIBUTES INCOMING APPLICATION TRAFFIC ACROSS MULTIPLE EC2 INSTANCES, IN MULTIPLE AVAILABILITY ZONES. THIS INCREASES THE FAULT TOLERANCE OF YOUR APPLICATIONS.
- ❖ THE LOAD BALANCER SERVES AS A SINGLE POINT OF CONTACT FOR CLIENTS, WHICH INCREASES THE AVAILABILITY OF YOUR APPLICATION. YOU CAN ADD AND REMOVE INSTANCES FROM YOUR LOAD BALANCER AS YOUR NEEDS CHANGE, WITHOUT DISRUPTING THE OVERALL FLOW OF REQUESTS TO YOUR APPLICATION. ELASTIC LOAD BALANCING SCALES YOUR LOAD BALANCER AS TRAFFIC TO YOUR APPLICATION CHANGES OVER TIME, AND CAN SCALE TO THE VAST MAJORITY OF WORKLOADS AUTOMATICALLY.
- ❖ YOU CAN CONFIGURE HEALTH CHECKS, WHICH ARE USED TO MONITOR THE HEALTH OF THE REGISTERED INSTANCES SO THAT THE LOAD BALANCER CAN SEND REQUESTS ONLY TO THE HEALTHY INSTANCES

TYPES OF LOAD BALANCERS

ELASTIC LOAD BALANCING SUPPORTS THREE TYPES OF LOAD BALANCERS:

- ❖ APPLICATION LOAD BALANCERS
- ❖ NETWORK LOAD BALANCERS
- ❖ CLASSIC LOAD BALANCERS

YOU CAN SELECT A LOAD BALANCER BASED ON YOUR APPLICATION NEEDS



LOAD BALANCER FLOW CHART

APPLICATION & CLASSIC LOAD BALANCER OVERVIEW

- AN APPLICATION LOAD BALANCER FUNCTIONS AT THE APPLICATION LAYER, THE SEVENTH LAYER OF THE OPEN SYSTEMS INTERCONNECTION (OSI) MODEL.
- AFTER THE LOAD BALANCER RECEIVES A REQUEST, IT EVALUATES THE LISTENER RULES IN PRIORITY ORDER TO DETERMINE WHICH RULE TO APPLY, AND THEN SELECTS A TARGET FROM THE TARGET GROUP FOR THE RULE ACTION USING THE ROUND ROBIN ROUTING ALGORITHM.
- WE CAN CONFIGURE LISTENER RULES TO ROUTE REQUESTS TO DIFFERENT TARGET GROUPS BASED ON THE CONTENT OF THE APPLICATION TRAFFIC. YOU CAN ADD ONE OR MORE LISTENERS TO YOUR LOAD BALANCER.
- YOU CAN CONFIGURE *HEALTH CHECKS*, WHICH ARE USED TO MONITOR THE HEALTH OF THE REGISTERED INSTANCES SO THAT THE LOAD BALANCER ONLY SENDS REQUESTS TO THE HEALTHY INSTANCES.

NETWORK LOAD BALANCER OVERVIEW

- A NETWORK LOAD BALANCER FUNCTIONS AT THE FOURTH LAYER OF THE OPEN SYSTEMS INTERCONNECTION (OSI) MODEL. IT CAN HANDLE MILLIONS OF REQUESTS PER SECOND.
- AFTER THE LOAD BALANCER RECEIVES A CONNECTION REQUEST, IT SELECTS A TARGET FROM THE TARGET GROUP FOR THE DEFAULT RULE.
- IT ATTEMPTS TO OPEN A TCP CONNECTION TO THE SELECTED TARGET ON THE PORT SPECIFIED IN THE LISTENER CONFIGURATION.

AZURE LOAD BALANCER

- AZURE LOAD BALANCER ALLOWS YOU TO SCALE YOUR APPLICATIONS AND CREATE HIGH AVAILABILITY FOR YOUR SERVICES. LOAD BALANCER SUPPORTS INBOUND AS WELL AS OUTBOUND SCENARIOS, AND PROVIDES LOW LATENCY, HIGH THROUGHPUT, AND SCALES UP TO MILLIONS OF FLOWS FOR ALL TCP AND UDP APPLICATIONS
- LOAD BALANCER WILL DISTRIBUTE NEW INBOUND FLOWS ARRIVING ON THE LOAD BALANCER'S FRONTEND TO BACKEND POOL INSTANCES ACCORDING TO RULES AND HEALTH PROBES.
- ADDITIONALLY, A PUBLIC LOAD BALANCER CAN ALSO PROVIDE OUTBOUND CONNECTIONS FOR VIRTUAL MACHINES INSIDE YOUR VIRTUAL NETWORK BY TRANSLATING THEIR PRIVATE IP ADDRESSES TO PUBLIC IP ADDRESSES.
- AZURE LOAD BALANCER IS AVAILABLE IN TWO DIFFERENT SKUS: BASIC AND STANDARD

STANDARD SKU VS BASIC SKU

	Standard SKU	Basic SKU
Backend pool size	up to 1000 instances	up to 100 instances
Backend pool endpoints	any virtual machine in a single virtual network, including blend of virtual machines, availability sets, virtual machine scale sets.	virtual machines in a single availability set or virtual machine scale set
Availability Zones	zone-redundant and zonal frontends for inbound and outbound, outbound flows mappings survive zone failure, cross-zone load balancing	/
Diagnostics	Azure Monitor, multi-dimensional metrics including byte and packet counters, health probe status, connection attempts (TCP SYN), outbound connection health (SNAT successful and failed flows), active data plane measurements	Azure Log Analytics for public Load Balancer only, SNAT exhaustion alert, backend pool health count
HA Ports	internal Load Balancer	/
Secure by default	default closed for public IP and Load Balancer endpoints and a network security group must be used to explicitly whitelist for traffic to flow	default open, network security group optional

STANDARD SKU VS BASIC SKU CONT....

Outbound connections	Multiple frontends with per rule opt-out. An outbound scenario <i>must</i> be explicitly created for the virtual machine to be able to use outbound connectivity. VNet Service Endpoints can be reached without outbound connectivity and do not count towards data processed. Any public IP addresses, including Azure PaaS services not available as VNet Service Endpoints, must be reached via outbound connectivity and count towards data processed. When only an internal Load Balancer is serving a virtual machine, outbound connections via default SNAT are not available. Outbound SNAT programming is transport protocol specific based on protocol of the inbound load balancing rule.	Single frontend, selected at random when multiple frontends are present. When only internal Load Balancer is serving a virtual machine, default SNAT is used.
Multiple frontends	Inbound and outbound	Inbound only
Management Operations	Most operations < 30 seconds	60-90+ seconds typical
SLA	99.99% for data path with two healthy virtual machines	Implicit in VM SLA
Pricing	Charged based on number of rules, data processed inbound or outbound associated with resource	No charge

GCP LOAD BALANCER

TYPES OF LOAD BALANCING:

1. GLOBAL EXTERNAL LOAD BALANCING

- ❖ HTTP(S) LOAD BALANCING
- ❖ SSL PROXY LOAD BALANCING
- ❖ TCP PROXY LOAD BALANCING

2. REGIONAL EXTERNAL LOAD BALANCING

- ❖ NETWORK LOAD BALANCING

3. REGIONAL INTERNAL LOAD BALANCING

- ❖ INTERNAL LOAD BALANCING

HTTP(S) LOAD BALANCING

- HTTP(S) LOAD BALANCING CAN BALANCE HTTP AND HTTPS TRAFFIC ACROSS MULTIPLE BACKEND INSTANCES, ACROSS MULTIPLE REGIONS.
- YOUR ENTIRE APP IS AVAILABLE VIA A SINGLE GLOBAL IP ADDRESS, RESULTING IN A SIMPLIFIED DNS SETUP.
- HTTP(S) LOAD BALANCING IS SCALABLE, FAULT-TOLERANT, REQUIRES NO PRE-WARMING, AND ENABLES CONTENT-BASED LOAD BALANCING. FOR HTTPS TRAFFIC, IT PROVIDES SSL TERMINATION AND LOAD BALANCING.
- HTTP(S) LOAD BALANCING SUPPORTS BOTH IPv4 AND IPv6 ADDRESSES FOR CLIENT TRAFFIC. CLIENT IPv6 REQUESTS ARE TERMINATED AT THE GLOBAL LOAD BALANCING LAYER, THEN PROXIED OVER IPv4 TO YOUR BACKENDS.
- HTTP REQUESTS CAN BE LOAD BALANCED BASED ON PORT 80 OR PORT 8080. HTTPS REQUESTS CAN BE LOAD BALANCED ON PORT 443

SSL LOAD BALANCING

- WITH SSL (TLS) PROXYING FOR YOUR SSL TRAFFIC, YOU CAN TERMINATE SSL SESSIONS AT THE GLOBAL LOAD BALANCING LAYER, THEN FORWARD THE TRAFFIC TO YOUR VIRTUAL MACHINE INSTANCES USING SSL (RECOMMENDED) OR TCP.
- SSL PROXY IS A GLOBAL LOAD BALANCING SERVICE. YOU CAN DEPLOY YOUR INSTANCES IN MULTIPLE REGIONS, AND THE LOAD BALANCER AUTOMATICALLY DIRECTS TRAFFIC TO THE CLOSEST REGION THAT HAS CAPACITY. IF THE CLOSEST REGION IS AT CAPACITY, THE LOAD BALANCER AUTOMATICALLY DIRECTS NEW CONNECTIONS TO ANOTHER REGION WITH AVAILABLE CAPACITY. EXISTING USER CONNECTIONS REMAIN IN THE CURRENT REGION.
- SSL PROXY PROVIDES SSL TERMINATION FOR YOUR NON-HTTPS TRAFFIC WITH LOAD BALANCING.
- WITH SSL (TLS) PROXYING FOR YOUR SSL TRAFFIC, YOU CAN TERMINATE SSL SESSIONS AT THE GLOBAL LOAD BALANCING LAYER, THEN FORWARD THE TRAFFIC TO YOUR VIRTUAL MACHINE INSTANCES USING SSL (RECOMMENDED) OR TCP

TCP PROXY LOAD BALANCING

- TCP LOAD BALANCING CAN SPREAD TCP TRAFFIC OVER A POOL OF INSTANCES WITHIN A COMPUTE ENGINE REGION.
- GOOGLE CLOUD PLATFORM (GCP) TCP PROXY LOAD BALANCING ALLOWS YOU TO USE A SINGLE IP ADDRESS FOR ALL USERS AROUND THE WORLD. GCP TCP PROXY LOAD BALANCING AUTOMATICALLY ROUTES TRAFFIC TO THE INSTANCES THAT ARE CLOSEST TO THE USER.
- CLOUD TCP PROXY LOAD BALANCING IS INTENDED FOR NON-HTTP TRAFFIC. FOR HTTP TRAFFIC, HTTP(S) LOAD BALANCING IS RECOMMENDED INSTEAD. FOR PROXIED SSL TRAFFIC, USE SSL PROXY LOAD BALANCING.
- TCP PROXY LOAD BALANCING SUPPORTS BOTH IPV4 AND IPV6 ADDRESSES FOR CLIENT TRAFFIC. CLIENT IPV6 REQUESTS ARE TERMINATED AT THE GLOBAL LOAD BALANCING LAYER, THEN PROXIED OVER IPV4 TO YOUR BACKENDS.
- IT IS SCALABLE, DOES NOT REQUIRE PRE-WARMING, AND HEALTH CHECKS HELP ENSURE ONLY HEALTHY INSTANCES RECEIVE TRAFFIC.

NETWORK LOAD BALANCING

- NETWORK LOAD BALANCING DISTRIBUTES TRAFFIC AMONG A POOL OF INSTANCES WITHIN A REGION. NETWORK LOAD BALANCING CAN BALANCE ANY KIND OF TCP/UDP TRAFFIC
- NETWORK LOAD BALANCING ALLOWS YOU TO BALANCE LOAD OF YOUR SYSTEMS BASED ON INCOMING IP PROTOCOL DATA, SUCH AS ADDRESS, PORT, AND PROTOCOL TYPE.
- NETWORK LOAD BALANCING USES FORWARDING RULES THAT POINT TO TARGET POOLS, WHICH LIST THE INSTANCES AVAILABLE FOR LOAD BALANCING AND DEFINE WHICH TYPE OF HEALTH CHECK THAT SHOULD BE PERFORMED ON THESE INSTANCES. SEE THE NETWORK LOAD BALANCING EXAMPLE FOR MORE INFORMATION.
- NETWORK LOAD BALANCING IS A REGIONAL, NON-PROXIED LOAD BALANCER. YOU CAN USE IT TO LOAD BALANCE UDP TRAFFIC, AND TCP AND SSL TRAFFIC ON PORTS THAT ARE NOT SUPPORTED BY THE SSL PROXY AND TCP PROXY LOAD BALancers.
- A NETWORK LOAD BALANCER IS A PASS-THROUGH LOAD BALANCER. IT DOES NOT PROXY CONNECTIONS FROM CLIENTS

INTERNAL LOAD BALANCING

- INTERNAL LOAD BALANCING DISTRIBUTES TRAFFIC FROM GOOGLE CLOUD PLATFORM VIRTUAL MACHINE INSTANCES TO A GROUP OF INSTANCES IN THE SAME REGION.
- INTERNAL LOAD BALANCING ENABLES YOU TO RUN AND SCALE YOUR SERVICES BEHIND A PRIVATE LOAD BALANCING IP ADDRESS WHICH IS ACCESSIBLE ONLY TO INSTANCES INTERNAL TO YOUR VIRTUAL PRIVATE CLOUD (VPC).
- INTERNAL LOAD BALANCING WORKS WITH AUTO MODE VPC NETWORKS, CUSTOM MODE VPC NETWORKS, AND LEGACY NETWORKS.
- INTERNAL LOAD BALANCING CAN ALSO BE IMPLEMENTED WITH REGIONAL MANAGED INSTANCE GROUPS. THIS ALLOWS YOU TO AUTOSCALE ACROSS A REGION, MAKING YOUR SERVICE IMMUNE TO ZONAL FAILURES.