Al-Based DDoS Attacks Judgment day 2024

f5

Ondrej Ciz F5, Solution Engineer III

Agenda

AI Evolution And Adoption

AI-Based DDoS Attack Architecture

AI-Based DDoS Attack Mitigation

Q&A



2 ©2024 F5

Al is not new

History of AI

1950 Turing Test	1955 The Term Al	1956 Perceptron K-Means Clusterir	ng	1960 GANs	1966 Eliza	Al Winte Genetic	– 1980s er Programming	1986 Backpropagatio
	• 2017 Automated recommenders surpass humans	2017 AlphaGo	2014 Alexa	2011 Watson	2006 Deep Neural Networks	2005 Self-Driving Vehicles	1997 Deep Blue LSTMs	1995 SVMs
	2018 OpenAl - GPT	2020 OpenAl launch	2(D/)21 ALL-E	2021 Cancer Identification	2022 ChatGP (GPT-3)	T 0 La V G	023 Open Source LLMs angChain Sector Databases

G

Al adoption set to define new business models, transform industries, global economies

80%

of tech executives will increase investment in Al in the next year¹ Al will likely contribute an estimated

\$15.7 Trillion

to the global economy by 2030³

The Europe Artificial Intelligence Market Size is Expected to reach

USD 325.29 Billion

12

by 2032⁵

AI will have an estimated

21%

net increase on United States GDP in 2030²

Asia's Top 1000 Companies to Allocate

Over 50%

of their IT Spending on AI Initiatives by 2025 according to IDC⁴

4 ©2024 F5

Top AI Challenges



Al is a business imperative but requires governance programs

Lack of visibility into AI use (Shadow AI)

Complexity of managing increasingly distributed AI workloads

Defending against AI attacks on apps

Safely connecting AI models with your data that lives in disparate environments

Protecting the explosion new API endpoints required to connect AI/ML models

SecOps' lack of visibility or governance of new API endpoints and middleware dependencies

Managing AI workload complexity without impacting business agility and innovation

Top business drivers and considerations for AI adoption



The Pace of Adoption of Generative Al Has Been Astounding

Time it took companies to reach 100 million users:



1

Sources: Global X ETFs with info derived from: BBC News. (2018, Jan 23). Netflix's history: From DVD rentals to streaming success; Cerullo, M. (2023, Feb 1). ChatGPT user base is growing faster than TikTok. CBS News.

7 ©2024 F5

DDoS attacks are on the rise significantly YoY 2022/2023

DDoS impacts include:

- Loss of revenue
- Loss of productivity (IT Ops / Security)
- Hiring specialized consultants
- Credits to consumers
- Legal and compliance fees
- Public relations

Indirect costs often include:

- Damage to the brand
- · Theft of vital data
- Loss of valuable customers and
- Opportunity loss

55 seconds between attacks in 202214 seconds between attacks in 2023



8 ©2023 F5 *https://www.link11.com/en/download/european-cyber-report-2023/

For Cyber Security AI Is a Double-Edged Sword

Creates new attack surfaces and new attack capabilities

Moving at alarming speed, no reservations



(F)

AI-Based DDoS Attack Architectures

Machine Learning-Enhanced Attackers: Training Models

- Understand and mimic normal network behavior, difficult to detect malicious activity.
- Adaptive Attack Strategies:
 - ML algorithms can adapt in real-time, allowing attackers to modify their strategies based on the effectiveness of the ongoing attack and the defensive measures in place.

• Neural Networks (GAN, VAE): Traffic Generation:

• GANs can be employed to generate realistic-looking traffic that can mimic legitimate user behavior, making it challenging for traditional rule-based defenses to distinguish between malicious and legitimate traffic.





Generative Adversarial Networks

Variational autoencoder

10 ©2023 F5 *<u>https://www.linkedin.com/pulse/future-ai-based-ddos-attacks-proactive-defense-pratik-barahatte-dahuf</u>

AI-Based DDoS Attack Architectures

Reinforcement Learning (RL) to optimize Attack Strategies

• Attackers might use RL algorithms to optimize the parameters of their attack in real-time, making it more difficult for defenders to predict and mitigate the attack.



Reinforcement Learning

IoT botnets

 IoT devices are often insufficiently protected and are therefore easily compromised. Compromised devices can be used to expand powerful botnet networks for DDoS attacks (e.g. Mirai botnet attacks of 2016). Al algorithms can be used to coordinate these botnets more efficiently and generate more sophisticated attack patterns.

11 ©2023 F5 *<u>https://www.linkedin.com/pulse/future-ai-based-ddos-attacks-proactive-defense-pratik-barahatte-dahuf</u>

Generative AI Has Democratized Technology



12 ©2024 F5

DDoS attack mitigation

F5 has been an AI pioneer for decades – and continues to evolve to meet market needs

2006 Policy Builder to automatically create WAF security policy

2011 Machine Learning (ML) to optimize BIG-IP network stack

2012 ML-based learning mode to improve WAF efficacy

2014

Behavioral anomaly analysis to improve WAF efficacy

2020

Shape acquisition Al powered Anti- Bot & Fraud

2019

Stage 2 retrospective analysis to deter bots and automated attacks

2016

Behavioral DoS to mitigate attacks using dynamic signatures

2015

ML in iHealth to identity customer BIG-IP issues and misconfigurations

2021

Dynamic API Discovery to mitigate risk from rogue APIs

2021

False Positive Suppression to reduce burden on SecOps teams

2022

Malicious User Detection to automatically mitigate bad actors

2023

API Endpoints Risk Score to measure risk of discovered APIs

Mitigation Strategies

Detection of Network Attacks using Machine Learning and Deep Learning Models



16

Mitigation Strategies

- Traffic Analysis and Anomaly Detection
- Behavioral Analysis
- Rate Limiting and Traffic Filtering
- Web Application Firewalls (WAF)
- Intrusion Detection and Prevention Systems (IDPS)
- Challenge-Response Mechanisms
- Cloud-Based DDoS Protection Services
- Regularly Update and Patch Systems
- Collaborate and Share Threat Intelligence
- Incident Response Planning



Volumetric DDoS Attacks



Protocol DDoS Attacks



Application DDoS attacks

The FLOW of Security [Logical Flow]



17 ©2023 F5

Auto DoS Protection

Enabled by default since Jan 2024 release

- Included as part of the base package.
- Detects anomalies for request rate, error rate, latency and throughput.
- Dynamically programming
 - Fast ACL to drop attacker at L3 layer
 - L7 ACL in envoy to drop connection based on fingerprint, geo or path
- Auto mitigation rules are created for a duration of 20 minutes, then auto deleted.
- Ability to create trusted client rules to bypass IP addresses from the mitigation.

DoS Protection

* L7 DDoS Auto Mitigation 🚯

℃ Default	
Default Default Block suspicious sources and serve JavaScript challenge to rest of traffic	
Block Block suspicious sources	
JavaScript Challenge Serve JavaScript challenge to suspicious sources	

Fast ACL rules

- Specified in terms of five tuple of the packet (dst ip, dst port, src ip, src port, protocol).
- Very efficient, but hard to do manually
- Block source IP attacker on L3
- Auto DoS Mitigation create them automatically

Fast ACL Type

* Select Site Type For acl 🚯

⊃‡ Site Type Regional Edge	^
Site Type Customer Edge ACL will be applied at customer edge sites	
Site Type Regional Edge ACL will be applied at regional edge sites	Default

(6

Order	Name	ACL Action	Source Ports	Source Prefixes	Ac
••• 1 ••• 1	block-source-ips	Deny		161.97.88.46/32 195.177.217.131/32 75.119.150.125/32 5.189.146.59/32 <u>19 more</u>	
			Add Item		

JavaScript Challenge

- Validate that the request is coming from a browser that is capable of running JavaScript
- Slows down a potential DoS attacks by forcing the browser to run a complex operation that requires many CPU cycles
- The load balancer tags response header with a cookie to avoid JavaScript challenge for subsequent requests.
- Can be configured via Policy Based Challenge rules for advance traffic matching criteria



	C v
Clients	
* Source IP Match 🚯	
ੋ੍ਹੈ Any Source IP	,
* Source ASN Match ③	
⊇‡ Any Source ASN	
* Client Selection (j)	
⊂; Any Client	
TLS Fingerprint Matcher ③	
▲ Not configured Configure >	

0

20 ©2023 F5

0 Be extremely careful with the JS challenge in XC. We turned it on and it caused a considerable amount of downtime for some customers. If Bot Protection is employed, do not use the JS challenge in XC. Also, never turn this on without conferring with the customer first. Just a few guidelines.

Alex Barajas; 2024-02-05T20:04:45.132

Rate Limiting

- Apply on IP Addresses by default, but can be configured for other identification criteria
- Optionally exempt rate limiting for requests from specified IP prefixes
- Trigger 429 codes in the logs
- The rate limit is always evaluated before any configured network security policy sets.

* Identifier Type 🚯	
Client IP Address	
Client IP Address Use client IP address as user identifier.	Default
Cookie Name Use the HTTP cookie value for the given name as user	identifier.
HTTP Header Name Use the HTTP header value for the given name as user	identifier.
Query Parameter Key Use the query parameter value for the given key as use	er identifier.
TLS Fingerprint Use TLS Fingerprint as user identifier.	
Client IP and HTTP Header Name Use the combination of Client IP and HTTP header valu given name as user identifier	ie for the

Rate Limit Configuration

^

* Nu	mber 🛈			
20	D			С×
* Pe	r Period	(j)		
Se	econd			×
5				
(s) A	llowed w	ithout Ra	ate Limiting 🚯	
(s) A IP	Allowed w Allowed	ithout Ra	ate Limiting 🚯	₽ C ∨
(s) A ≱ IP	Allowed w Allowed	ithout Ra List Order	ate Limiting ③	€ C ∨
(s) A ≱ IP	Allowed w Allowed	ithout Ra List Order 1	IPv4 Prefix List ③	₽ C ∨ ×

IP Reputation

Database of known malicious IP addresses classified by threat categories

- Spam sources
- Mobile threats
- Windows exploits
- Web attacks
- Botnets
- Scanners
- Denial of Service (DoS)
- Phishing

Form Documentation JSON		hightarrow Reset All Fields Q s
Edit HTTP Load Balancer frontend	Common Security Controls	
Metadata Domains and LB Type	* Service Policies (2) and Apply Specified Service Policies × ~	
Origins Routes Web Application Firewall	 *Apply Specified Service Policies ⊙ ○ Configured Edit Configuration > 	
API Protection DoS Protection Client-Side Defense Common Security Controls	IP Reputation ③ [®] → [®] Enable × · * List of IP Threat Categories to choose ③	
Other Settings	Spam Sources Phishing X V Please fill out this field.	
	* User Identifier ③ \$\star\$_\$^a Client IP Address ×	
	* Malicious User Detection ③ Implie × v	
Cancel and Exit	* Rate Limiting ①	

22 ©2023 F5

F5 Distributed Cloud DDoS Mitigation Network

Responds to DDoS attacks in < 2 minutes on average.* Top BGP peering

Modern Service Provider Global Backbone

Flexible Service Options including Always Available or Always On deployments

Connect how and where you need with BGP or Proxy-based traffic redirection and direct connections, peering or GRE tunnels for clean traffic return.



Standard DDoS Service offering MSA specifies a 15 Minute Response SLA.

F

Metrics request rate & spikes



24 ©2023 F5

DDoS Events map

- Shows details on DDoS events occurring over selected time interval.
- Display DDoS trend including start time, end time, and number of events represented.



F

Why F5 for API Security?



Discovery + Enforcement

Holistic app and API security, with AI/ML powered continuous discovery to uncover unknown, forgotten and shadow APIs paired with critical protection capabilities (WAF, DoS mitigation, Bot protection etc.)



Flexibility and Architectural freedom

Combination of hardware, software and SaaS with flexible, lightweight, agnostic tools customers can deploy across cloud, on-premises, and edge environments.



Goes Beyond Detection and Blocking

A portfolio that can solve the problem more holistically giving Platform Ops and IT teams the tools to enable developers with self-service capabilities to discover, use, publish and manage APIs and security teams consistent oversight and fine-grained controls over all apps and API endpoints.





EXTRA SLIDES

(F)

29 ©2024 F5 CONFIDENTIAL

Generative AI: The Democratization of AI

From tool for a few, to resource for many

Generative AI takes the main stage



31 ©2024 F5 CONFIDENTIAL

G

GenAl APIs/GenAl-enabled apps in use already with widely accelerated growth expected over the next few years



More than half of tech executives whose companies are experimenting with generative AI (56%) are doing so for economic savings¹ pwc

Nearly all business leaders say their company is prioritizing at least one initiative related to AI systems in the near term, while 54% of companies have already implemented GenAI in some areas of their business² Gartner

Press Release STAMFORD, CT, October 11, 2023

By 2026, more than 80% of enterprises will have used generative artificial intelligence (GenAI) application programming interfaces (APIs) or models, and/or deployed GenAIenabled applications in production environments, up from less than 5% in 2023.³

32 ©2024 F5 CONFIDENTIAL

Generative AI Has Engendered both Excitement and Caution



"AI will be the biggest thing this decade...it has potential to teach math and offer medical advice to people with limited access to resources."

Bill Gates, Cofounder of Microsoft



"Humans must be unambiguously, unquestionably in charge of powerful AI models to prevent them from going out of control."

Satya Nadella, CEO of Microsoft



"We see enormous potential in this space to affect virtually everything we do... It will affect every product and every service that we have." **Tim Cook**, CEO at Apple



"Al can have huge problems, in terms of a democracy and how it reacts to this... biggest value is in replacing human labor."

Warren Buffet, Business magnate, investor, and philanthropist

Generative Al Multicloud Networking Challenges

Connecting Customer Edge nodes directly to each other for internal apps with secure service connectivity, or via the Regional Edge for publicfacing apps that need a front door service.

Connecting and sharing data between inference apps, the central LLM, and training data through secure network connectivity.

Enablng inference app deployment and maintenance at scale with an ingress controller in each app cluster.

Enabling load balancing and security at each node, with centralized DDoS, Bot, and API protection, plus distributed web app firewalls locally.











Challenges with Running AI/ML Workloads in Kubernetes

Across hybrid, multi-cloud environments with disaggregated technologies

- Connection timeouts and errors
- Insufficient visibility into app health and performance
- Difficulties with securing distributed app environments
- Increasing complexity and tool sprawl





 Poor user experiences

• Troubleshooting difficulties and downtime

- Increased risk of exposure to cybersecurity threats
- Hard to operate, manage, and troubleshoot







The challenge of connecting AI inference apps

Traditionally, LLMs, training data, and inference apps are all deployed centrally, and any user accessing them is routed to a single, centralized location.

Now, those inference apps are deployed in distributed environments, and are decoupled from the LLM.

However, **secure connectivity is still required** to protect the LLM, training data, and apps.



LLM, Training Data, & Inference App 0 Public Cloud A Inference App Inference App Data Center / Private Cloud Public Cloud / Edge B, C, D. **F5 Distributed Cloud** Secure Multi-Cloud Networking F5 Distributed Cloud Web App and API Protection T

Enabling connectivity across AI inference apps

Secure Multi-Cloud Networking ensures that inference apps can be deployed quickly to any cloud, edge, or on-prem site.

Each app has network connectivity already provisioned, with consistent L7 security, and visibility into traffic flows and threats across the network.



F5 Distributed Cloud Secure Multi-Cloud Networking App Connect

Discover > Deliver > Secure > Operate Any Inference App across Any Environment

Connect AI apps to speed up deployment at the edge





Key Capabilities for AI Apps

Secure, flexible deployment of inference opps, with connectivity to LLMs and data in any public cloud, on-premises, or edge site. App Connect offers:

- Service discovery
- App segmentation
- Service networking
- Multi-cluster ingress, egress, and deployment
- End-to-end encryption
- End-to-end observability
- Standardized app security integration across sites (e.g. WAAP, Bot Defense)

6

0 Same comment as slide 30 (this looks like key capabilities). Byron McNaught; 2024-01-12T19:55:59.799

F5 Distributed Cloud Secure Multi-cloud Networking Network Connect

Connect > Segment > Secure > Operate Networks across Any Environment

Connect Networks to secure AI Apps





Key Capabilities for AI Apps

Easily and securely network across public clouds, hybrid clouds, and edge sites, wherever your A and LLM workloads are running. Network Connect offers:

- Cloud orchestration and private connectivity
- Centralized resource access for identity control and shared data
- Global network segmentation
- Embedded network security and firewall service insertion at CEs
- End-to-end observability
- App networking integration

Snímka 39

Should this be "key capabilities for AI Apps"? Byron McNaught; 2024-01-12T19:55:38.902 0

00 Yeah, that works Colin Clauset; 2024-01-18T16:14:45.549

Generative Al Multi-Cloud Networking Challenges

- Connecting Customer Edge nodes directly to each other for internal apps with secure service connectivity, or via the Regional Edge for public-facing apps that need a front door service.
- Connecting and sharing data between inference apps, the central LLM, and training data through secure network connectivity.
- · Enabling inference app deployment and maintenance at scale with an ingress controller in each app cluster.
- · Enabling load balancing and security at each node, with centralized DDoS, Bot, and API protection, plus distributed web app firewalls locally.







Inference App Cluster (App, Inference Model, Ingress Controller)







Customer Edge

40 ©2024 F5

LLM



6

The challenge of connecting Al inference apps

Traditionally, LLMs, training data, and inference apps are all deployed centrally, and any user accessing them is routed to a single, centralized location.

Now, those inference apps are deployed in distributed environments, and are decoupled from the LLM.

However, secure connectivity is still required to protect the LLM, training data, and apps.

Training Data

Secure Network &

Service Connectivity

Customer

Edge

CE



6

41 ©2024 F5

LLM

Web App &

API Protection

Key

Enabling connectivity across AI inference apps

Secure Multicloud Networking

ensures that inference apps can be deployed quickly to any cloud, edge, or on-prem site.

Each app has network connectivity already provisioned, with consistent L7 security, and visibility into traffic flows and threats across the network.



Robust App Security Click-to-integrate full stack of app security controls

Flexible Network Security Network-layer security with support for third-party firewalls

End-to-End Private Network

End-to-end encryption across the F5 global network

Intent-based Policy Enforcement

Easily define, deploy, and enforce a single intent-based security policy across all apps and environments

Rich Observability

Network and application layer performance and security metrics

F5 Distributed Cloud Secure Multicloud Networking Network Connect

Connect > Segment > Secure > Operate Networks across Any Environment



Connect Networks to secure AI Apps

Key Capabilities for AI Apps

Easily and securely network across public clouds, hybrid clouds, and edge sites, wherever your AI and LLM workloads are running. Network Connect offers:

- Cloud orchestration and private connectivity
- Centralized resource access for identity control and shared data
- Global network segmentation
- Embedded network security and firewall service insertion at CEs
- End-to-end observability
- App networking integration

F5 Distributed Cloud Secure Multicloud Networking **App Connect**

Discover > Deliver > Secure > Operate Any Inference App across Any Environment



44 ©2024 F5

Customer Edge Use-Case : AI/ML-as-a-service Challenges



Customer Edge Use-Case : AI/ML-as-a-service



