

API Connect and DataPower Gateway

Performance Best
Practices



IBM Cloud

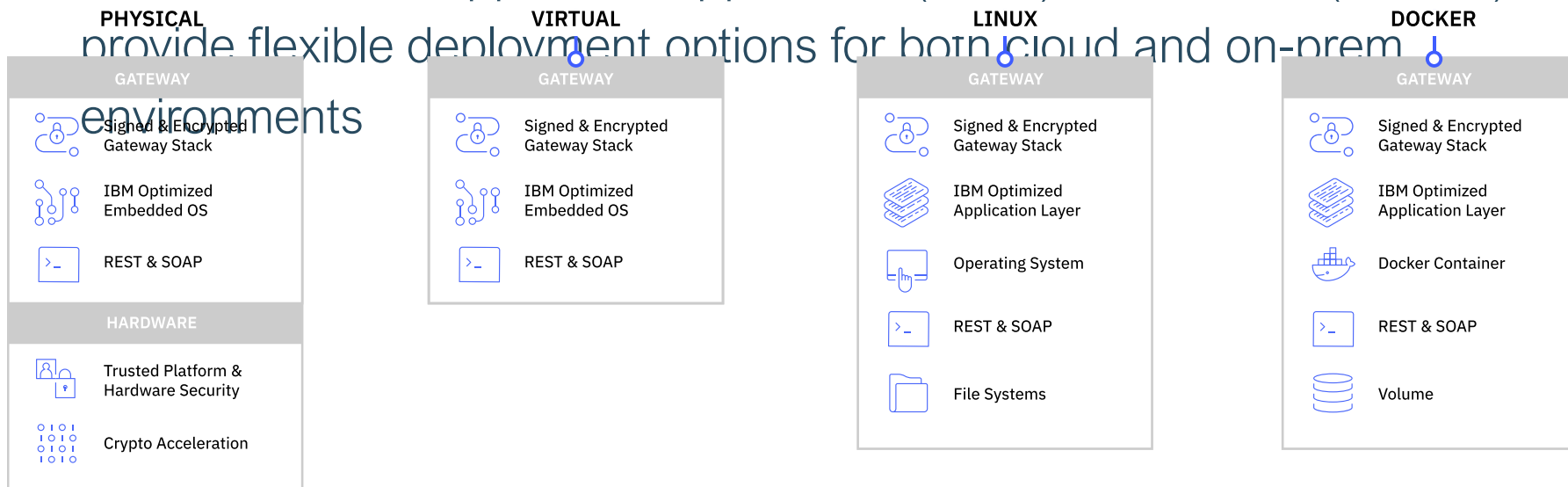


Please note

- IBM's statements regarding its plans, directions, and intent are subject to change or withdrawal without notice and at IBM's sole discretion.
- Information regarding potential future products is intended to outline our general product direction and it should not be relied on in making a purchasing decision.
- The information mentioned regarding potential future products is not a commitment, promise, or legal obligation to deliver any material, code or functionality. Information about potential future products may not be incorporated into any contract.
- The development, release, and timing of any future features or functionality described for our products remains at our sole discretion.
- Performance is based on measurements and projections using standard IBM benchmarks in a controlled environment. The actual throughput or performance that any user will experience will vary depending upon many factors, including considerations such as the amount of multiprogramming in the user's job stream, the I/O configuration, the storage configuration, and the workload processed. Therefore, no assurance can be given that an individual user will achieve results similar to those stated here.

DataPower Gateways can deploy anywhere...

- **Physical appliances:** All-in-one (HW / SW), DMZ-ready with physical security including crypto acceleration and optional hardware security module (HSM)
- **Software:** virtual appliance, application (Linux) & container (Docker)



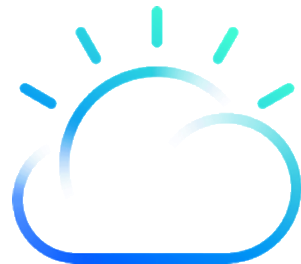
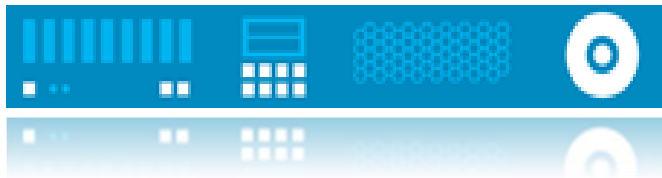
Choosing the right Gateway form factor

Physical appliances provides the most comprehensive security combined physical with firmware security.

Virtual, Linux and Container offer “right sized” units of capacity, as few as 4 CPUs

Container provide ability to leverage auto-scaling and runtime health monitoring

Container is “cloud ready” to facilitate both public and private cloud based deployments



Single API Gateway supports 30K TPS with 8 ms latency!

Natively built API Gateway using purpose-built technology for native OpenAPI/Swagger REST and SOAP APIs

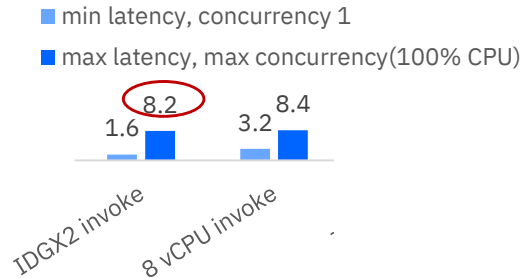
Multi-cloud scalability and extensibility to help meet SLAs and improve client user experience

IDG X2 physical appliances use the equivalent of 48 vCPU

Max Throughput @ 100% CPU

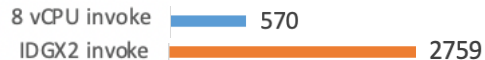


Latency(mS) @ Min/Max Concurrency



10X increased performance with natively built API Gateway

Before: DP Multi protocol
Gateway Service

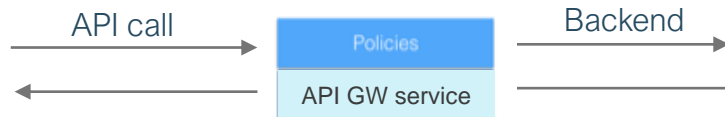


10X increased performance with natively built API
Gateway for both virtual and physical appliances

New: Native API
Gateway Service



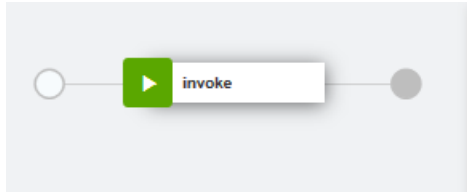
Throughput (TPS)



Benchmark APIs

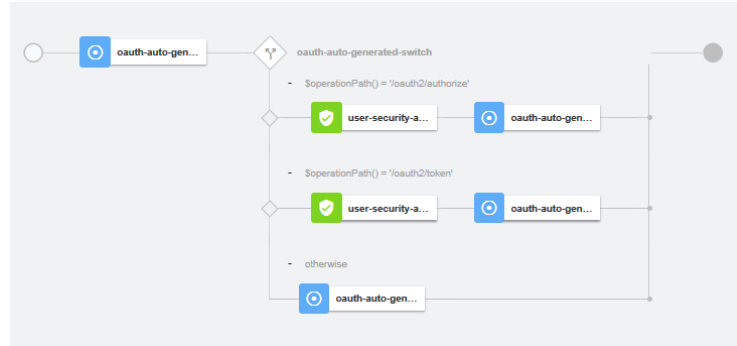
Invoke (Simple)

Invoke policy



Invoke + Security (Medium)

Invoke Policy +
OAuth Security



Invoke + Transformation (Complex)

Map +
Invoke Policy +
Map



Each Benchmark API Assembly includes:

- Rate-limit
- Activity-log
- Client-id validation

Each Benchmark APIs runs with

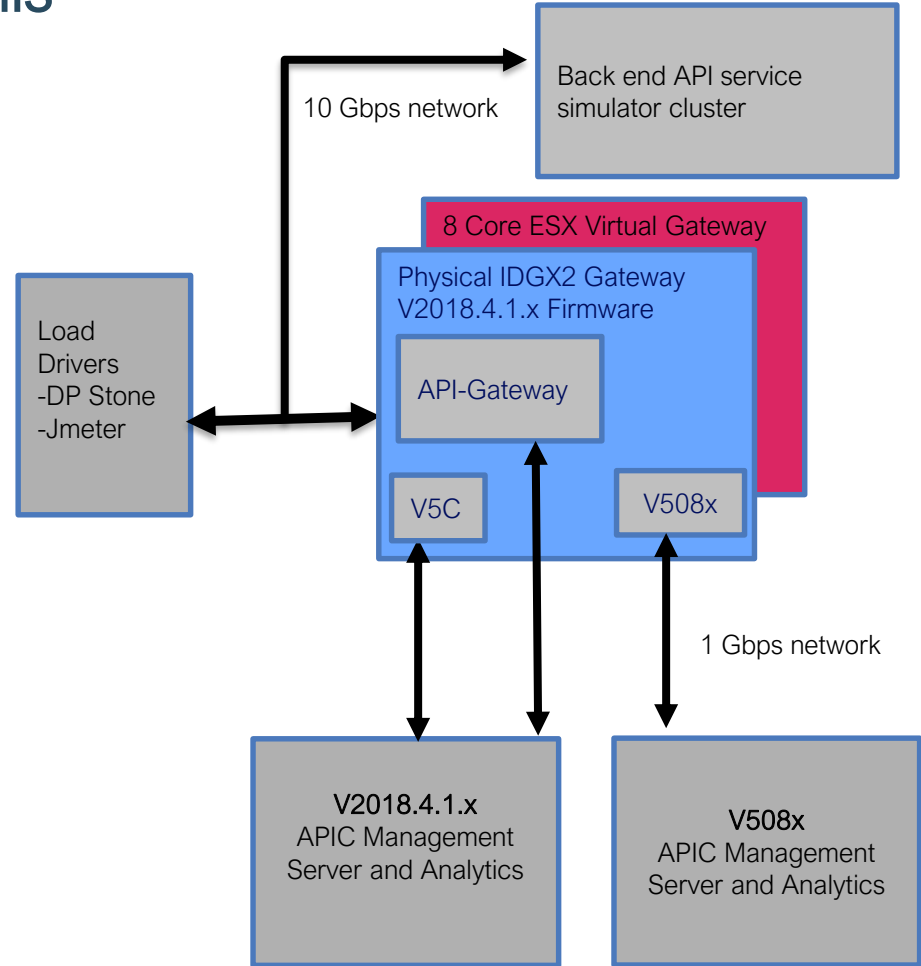
- 4KB request payload
- Running on IDGX2 and 8 vCPU virtual gateways

Performance Measurement Details

Single gateway service used for comparative analysis across three gateway types:

- API Gateway
- V2018 Multi-protocol Gateway (v5c)
- V5 Multi-protocol Gateway

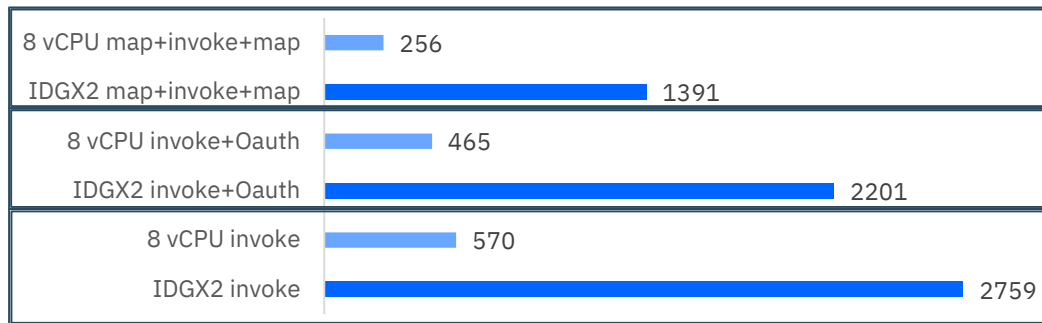
Backend service simulator is a “no-op” that delays response using random uniform distribution between 45 mS and 55 mS



API Connect V5 with v2018 MPGW Gateway

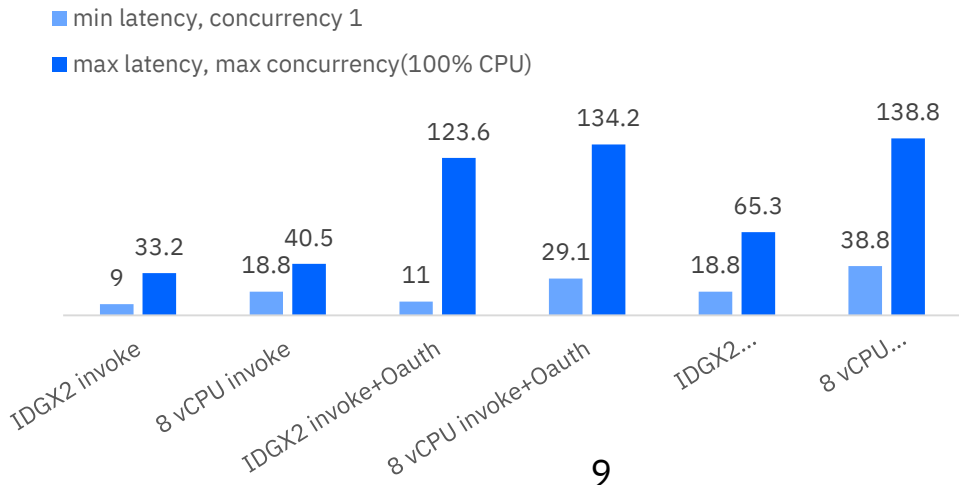
API Connect v5084 running with DataPower version 2018.4.1.1

Max Throughput @ 100% CPU



Throughput (TPS)

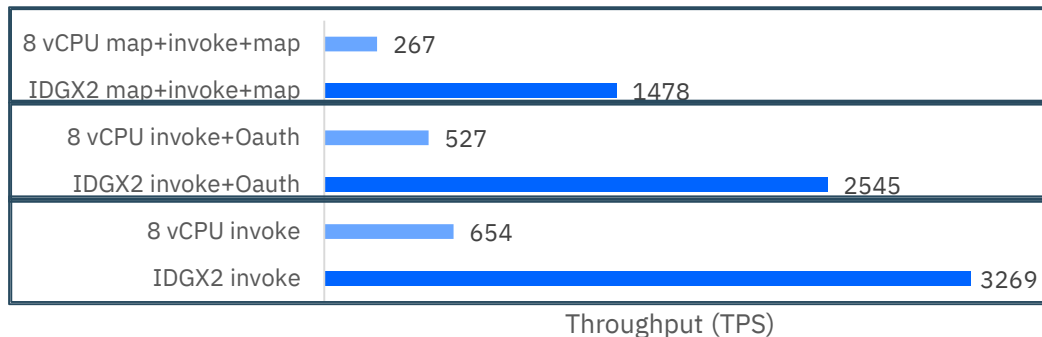
Latency(mS) @ Min/Max Concurrency



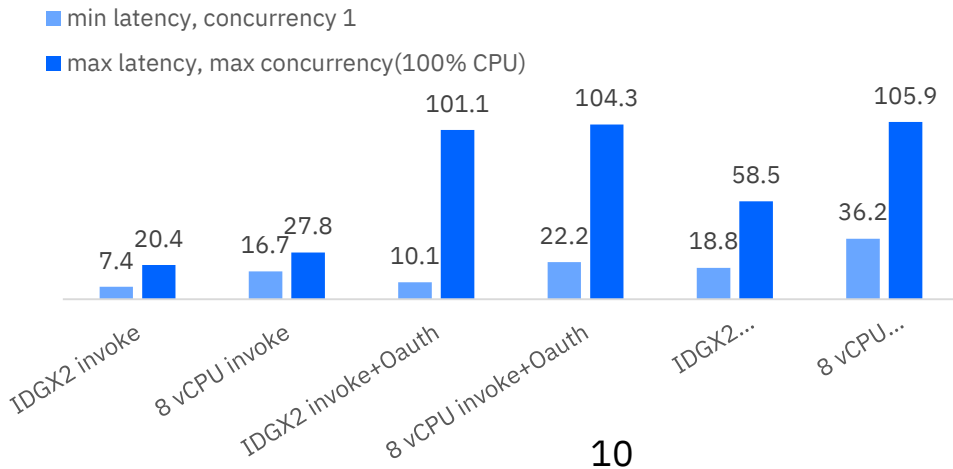
API Connect V2018 with v2018 MPGW Gateway

API Connect v2018.4.1 with DataPower version 2018.4.1.1, running V5c (Multi-protocol Gateway)

Max Throughput @ 100% CPU



Latency(mS) @ Min/Max Concurrency

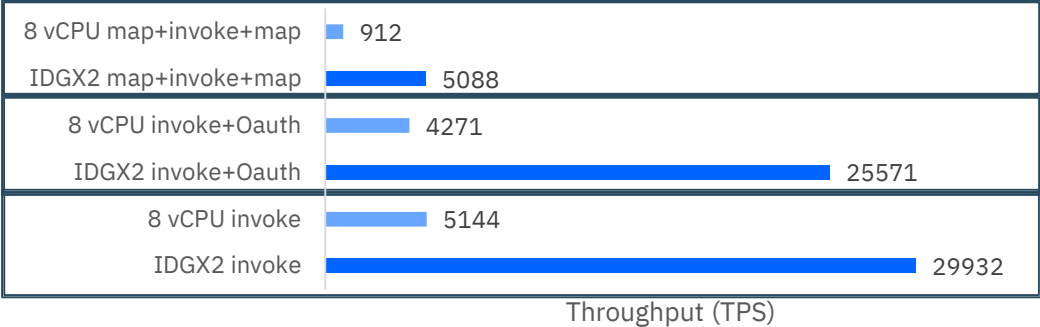


API Connect V2018 with v2018 API Gateway Service

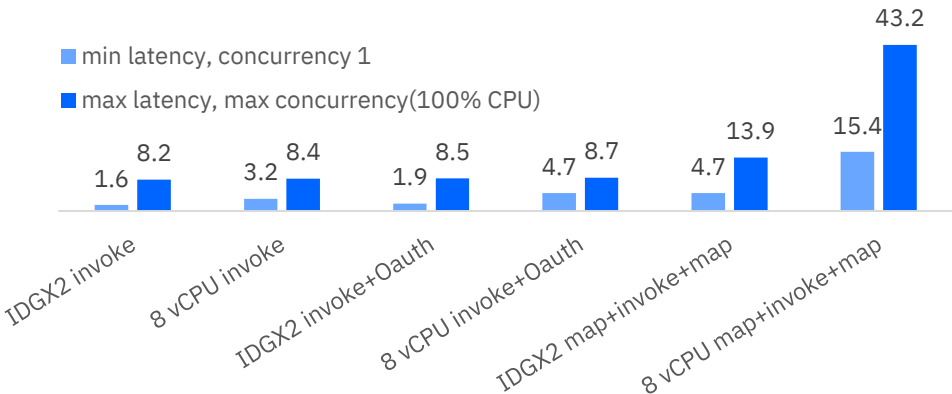
API Connect v2018.4.1 with DataPower version 2018.4.1.1, running API Gateway Service

Container form factor has same throughput and latency characteristics as Virtual using same underlying hardware

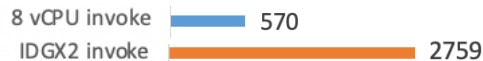
Max Throughput @ 100% CPU



Latency(mS) @ Min/Max Concurrency



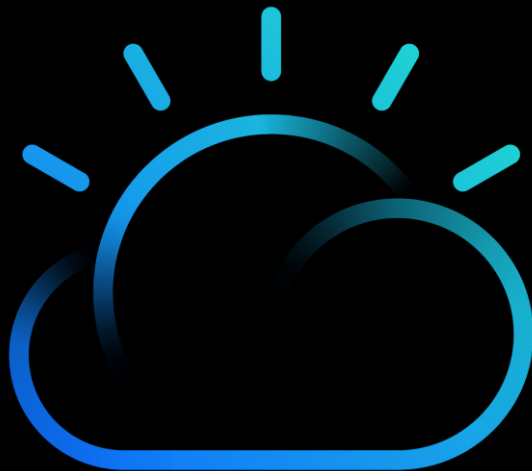
V2018 API Gateway Service vs V5 Multi-Protocol Gateway Service



10X increased performance with natively built API Gateway for both virtual and physical appliances



Thank You



Reserve Capacity For Planned and Unplanned Outages

HA events include planned and unplanned events

Best practices recommends HA planning with an assumption of 2 concurrent members failed, one planned and one unplanned

Active-active solutions are required

Each member of the cluster MUST have a reserve capacity in order to absorb the workload no longer being handled by failed members

Under an assumption that A% resource utilization is the upper limit under cluster member failure conditions

The equation to calculate resource utilization ceiling, during non failure conditions, as a function of number of cluster members N, with B members failed, is:

$$A = \left(\frac{n}{n-B} \right) * x$$

where: A is the maximum utilization by surviving members during a failure event, typically 80%

n is the number of cluster members

B is the number of failed cluster members

x is the maximum utilization in absence of member failure(s)

Substituting A=.8, B=1, and solving for x leaves:

$$x = \frac{.8n - .8}{n}$$

Substituting A=.8, B=2, and solving for x leaves:

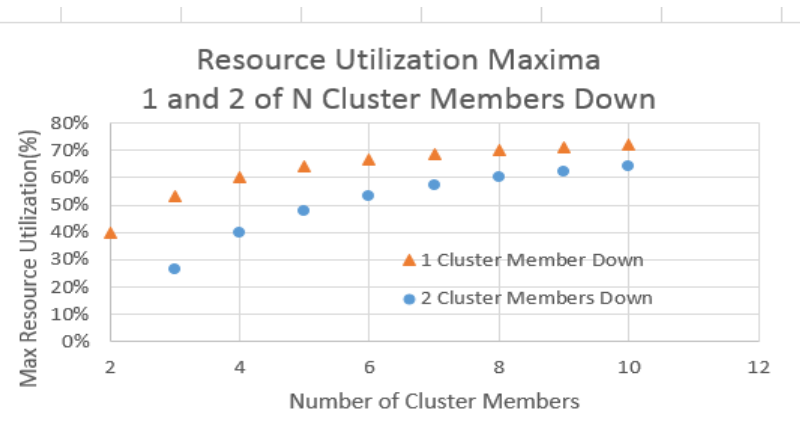
$$x = \frac{.8n - 1.6}{n}$$

Reserve Capacity For Planned and Unplanned Outages

Single member failure case:

- 2 member HA solution there is a 40% resource utilization ceiling, 60% reserve
Therefore, when one member fails the other member will continue at 80% utilization
- 3 member solution we have a 53% ceiling and when 1 member fails its 53% is divided evenly on the 2 remaining members which would then run at 79.5%
- 10 member solution has 72% ceiling and we have $72 + 72 / 9 = 80\%$

Max utilization during HA failure event:		0.8
Number of nodes down:	1	2
Number of nodes		
2	40%	
3	53%	27%
4	60%	40%
5	64%	48%
6	67%	53%
7	69%	57%
8	70%	60%
9	71%	62%
10	72%	64%

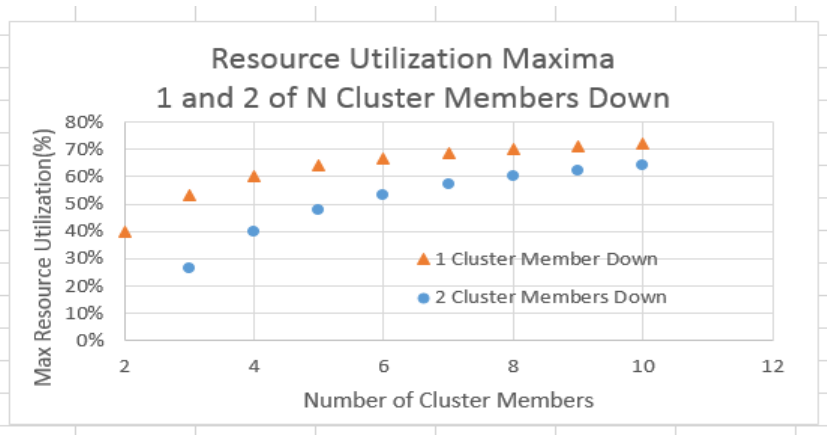


Reserve Capacity For Planned and Unplanned Outages

Dual member failure case:

- First, it is obvious that for the dual failure case we have to start with minimum of 3 members
- 3 member solution with 2 failures has a 27% ceiling, so the utilization on the single remaining member CPU is $27 + (2 \times 27) = 81\%$
- 4 member solution has ceiling of 40%, so the utilization on the 2 remaining members CPU is $40 + ((2 \times 40) / 2) = 80\%$
- 10 member solution, 8 remaining members CPU is $64 + ((2 \times 64) / 8) = 80\%$

Max utilization during HA failure event:		0.8
Number of nodes down:	1	2
Number of nodes		
2	40%	
3	53%	27%
4	60%	40%
5	64%	48%
6	67%	53%
7	69%	57%
8	70%	60%
9	71%	62%
10	72%	64%



Little's Law

Relationship Between Throughput and Latency

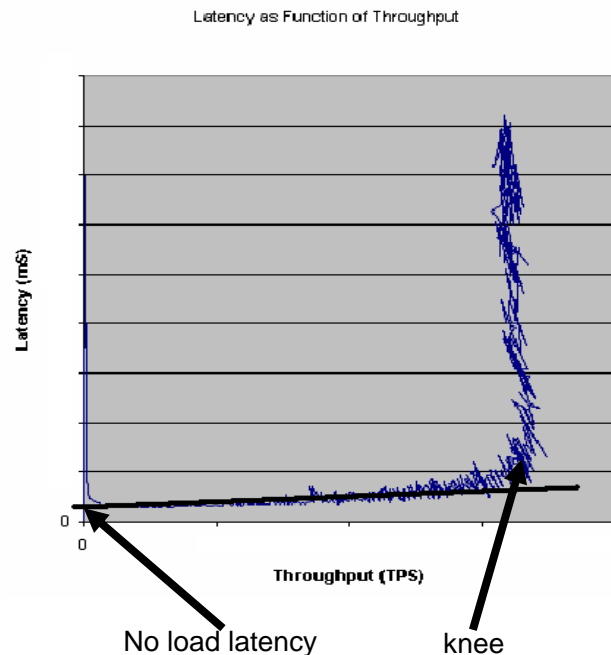
This characteristic pattern was proven as mathematical law by Professor John Little of MIT. More on Little's Law in next slide

It is helpful to describe a few key points on this characteristic pattern

- The “no load” latency:
This is the latency when throughput is minimal. By definition it can not be 0. It is the Y-Intercept.
- The “linear region” of system operation:
This is the portion of the graph where the throughput increases linearly with concurrency
- The “knee” value of TPS:
This is the maximum value of throughput where added concurrency only adds latency. The peak value of measured TPS is the saturation value.

Solution architects work to push the knee as far to the right as required by providing adequate infrastructure: CPU count, disk IO rate, network bandwidth and memory

Product architects and development teams work to move the knee as far to the right as possible for a given set of CPU/disk/network/memory resources, this also tends to reduce the slope and Y intercept of the linear region



Little's Law

Relationship Between Throughput and Latency

Little's Law: (https://en.wikipedia.org/wiki/Little's_law)

occupancy = latency x throughput

Applied to transactional systems:

concurrency = latency x throughput
= $S \times T/S$

The concurrency control in the load drivers, Jmeter for example, is generally described as number of “virtual users”

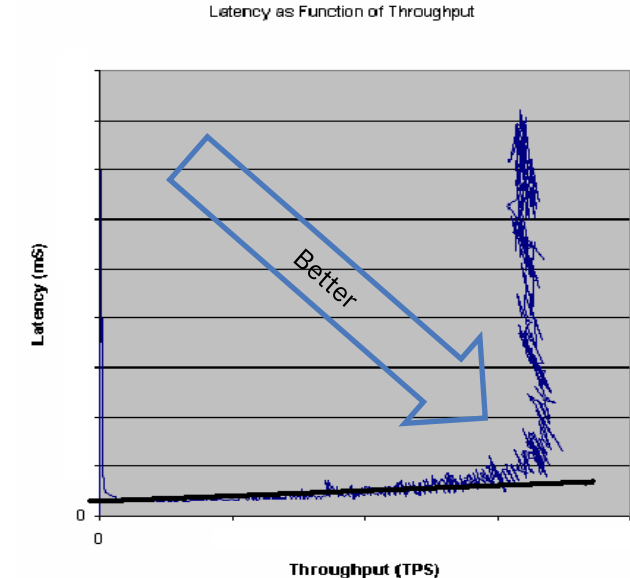
- Each virtual user has 1 and only 1 transaction live in the system at any time
- Nth transaction will not initialize until the N-1 transaction completes

When the system is stable the arrival rate = the completion rate

- The concurrency control and transaction latency govern the throughput
- The greater number of users, the greater the transaction arrival rate, the greater the completion rate

When the system is unstable, the throughput can not increase due to a constraint, then the arrival rate > completion rate

- 100% CPU, or network link speed, or IOPS capacity of disks, or memory utilization
- add more users, increased transaction arrival rate, the TPS can not increase so the latency MUST increase and the transaction completion rate saturates at constant a value

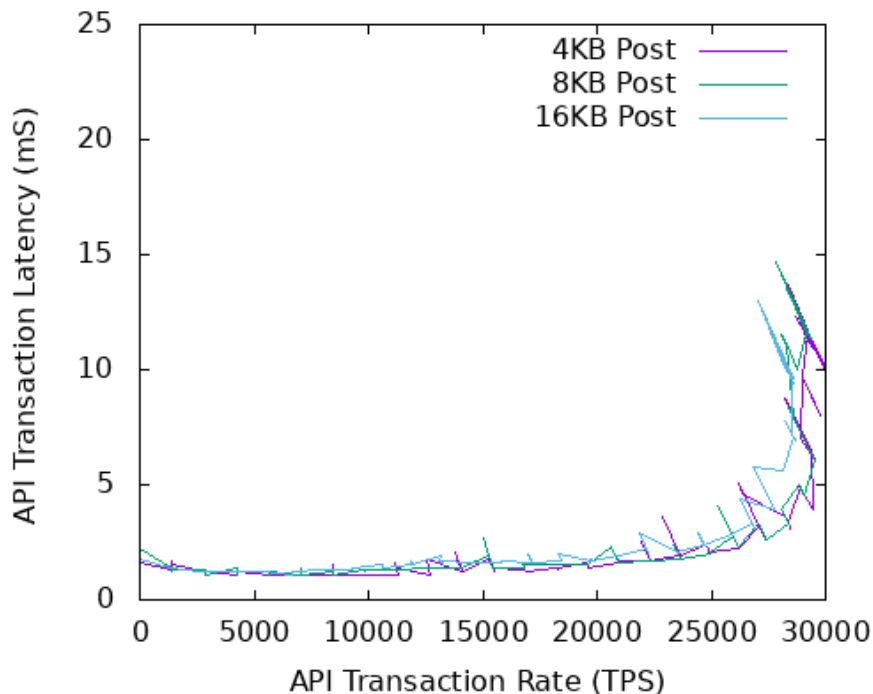


APIC Runtime on DP view of Little's Law

API Gateway Mode

- Note very small latency growth as TPS increases in the linear region of operation, CPU < 80%
- “no-load” latency ~2 mS
- “knee” throughput of 30K TPS
- “knee” latency of < 10 mS

APIC 2018.4.1.1-GA Latency as function of Throughput
IDGX2 Physical Platform
Invoke Assembly APIGW Mode
50mS backend latency



Notices and disclaimers

© 2018 International Business Machines Corporation. No part of this document may be reproduced or transmitted in any form without written permission from IBM.

U.S. Government Users Restricted Rights — use, duplication or disclosure restricted by GSA ADP Schedule Contract with IBM.

Information in these presentations (including information relating to products that have not yet been announced by IBM) has been reviewed for accuracy as of the date of initial publication and could include unintentional technical or typographical errors. IBM shall have no responsibility to update this information. **This document is distributed “as is” without any warranty, either express or implied. In no event, shall IBM be liable for any damage arising from the use of this information, including but not limited to, loss of data, business interruption, loss of profit or loss of opportunity.** IBM products and services are warranted per the terms and conditions of the agreements under which they are provided.

IBM products are manufactured from new parts or new and used parts. In some cases, a product may not be new and may have been previously installed. Regardless, our warranty terms apply.”

Any statements regarding IBM's future direction, intent or product plans are subject to change or withdrawal without notice.

Performance data contained herein was generally obtained in a controlled, isolated environments. Customer examples are presented as illustrations of how those customers have used IBM products and the results they may have achieved. Actual performance, cost, savings or other results in other operating environments may vary.

References in this document to IBM products, programs, or services does not imply that IBM intends to make such products, programs or services available in all countries in which IBM operates or does business.

Workshops, sessions and associated materials may have been prepared by independent session speakers, and do not necessarily reflect the views of IBM. All materials and discussions are provided for informational purposes only, and are neither intended to, nor shall constitute legal or other guidance or advice to any individual participant or their specific situation.

It is the customer’s responsibility to insure its own compliance with legal requirements and to obtain advice of competent legal counsel as to the identification and interpretation of any relevant laws and regulatory requirements that may affect the customer’s business and any actions the customer may need to take to comply with such laws. IBM does not provide legal advice or represent or warrant that its services or products will ensure that the customer follows any law.

Notices and disclaimers continued

Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements or other publicly available sources. IBM has not tested those products about this publication and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products. IBM does not warrant the quality of any third-party products, or the ability of any such third-party products to interoperate with IBM's products. **IBM expressly disclaims all warranties, expressed or implied, including but not limited to, the implied warranties of merchantability and fitness for a purpose.**

The provision of the information contained herein is not intended to, and does not, grant any right or license under any IBM patents, copyrights, trademarks or other intellectual property right.

IBM, the IBM logo, ibm.com and [names of other referenced IBM products and services used in the presentation] are trademarks of International Business Machines Corporation, registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on the Web at "Copyright and trademark information" at: www.ibm.com/legal/copytrade.shtml.