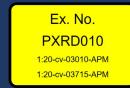


Prof. James Mickens

Distributed Systems and Security Expert



Assignment

IN THE UNITED STATES DISTRICT COURT FOR THE DISTRICT OF COLUMBIA

UNITED STATES OF AMERICA, et al.,

Plaintiffs,

Case No. 1:20-cv-03010-APM HON. AMIT P. MEHTA

Case No. 1:20-cy-03715-APM

GOOGLE LLC,

Defendant.

pt-i-

STATE OF COLORADO, et al.,

GOOGLE LLC,

HON, AMIT P. MEHTA

Defendant.

PLAINTIFFS' REVISED PROPOSED FINAL JUDGMENT

WHEREAS, Plaintiffs United States of America, and the States and Commonwealths of Arkansas, California, Georgia, Florida, Indiana, Kentucky, Louisiana, Michigan, Missouri, Mississippi, Montana, South Carolina, Texas, and Wisconsin, by and through their respective Attorneys General ("Co-Plaintiff States"), filed their Complaint on October 20, 2020, and their Amended Complaint on January 15, 2021;

AND WHEREAS, Plaintiffs Colorado, Nebraska, Arizona, Iowa, New York, North
Carolina, Tennessee, Utah, Alaska, Connecticut, Delaware, District of Columbia, Guam, Hawaii,
Idaho, Illinois, Kansas, Maine, Maryland, Massachusetts, Minnesota, Nevada, New Hampshire,
New Jersey, New Mexico, North Dakota, Ohio, Oklahoma, Oregon, Pennsylvania, Puerto Rico,

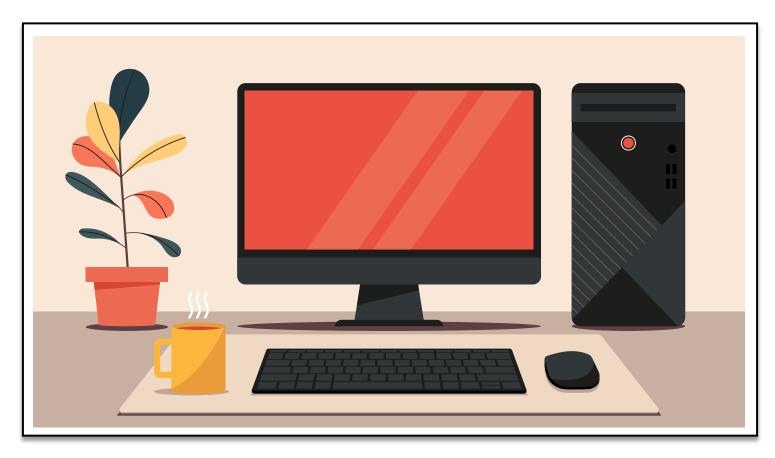
Assignment 1: Assess <u>divesting Chrome</u> from a technical perspective

Assignment 2: Assess <u>data sharing and</u> <u>syndication</u> remedies from a technical perspective

Why Distributed Systems?

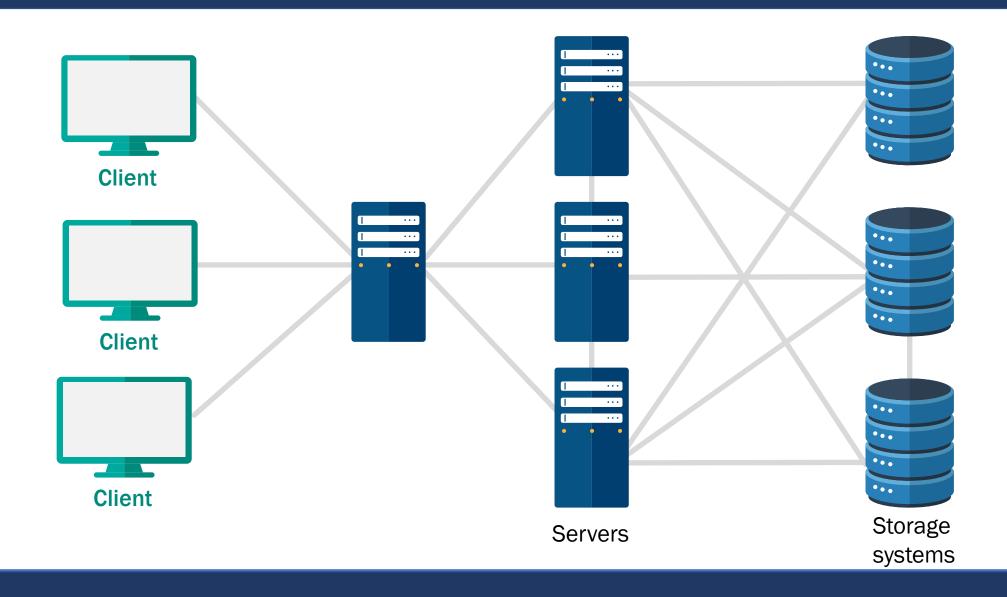
Understanding distributed systems is a necessary step in understanding the divestiture, data sharing, and syndication remedies.

Distributed Systems: An Overview



Single Computer

Distributed Systems: An Overview



Key Conclusions

Conclusion 1: From a technical perspective, Chrome divestiture is feasible.

Conclusion 2: From a technical perspective, data sharing and syndication remedies are feasible.

Technical Feasibility



Type of work Google and third parties already do.

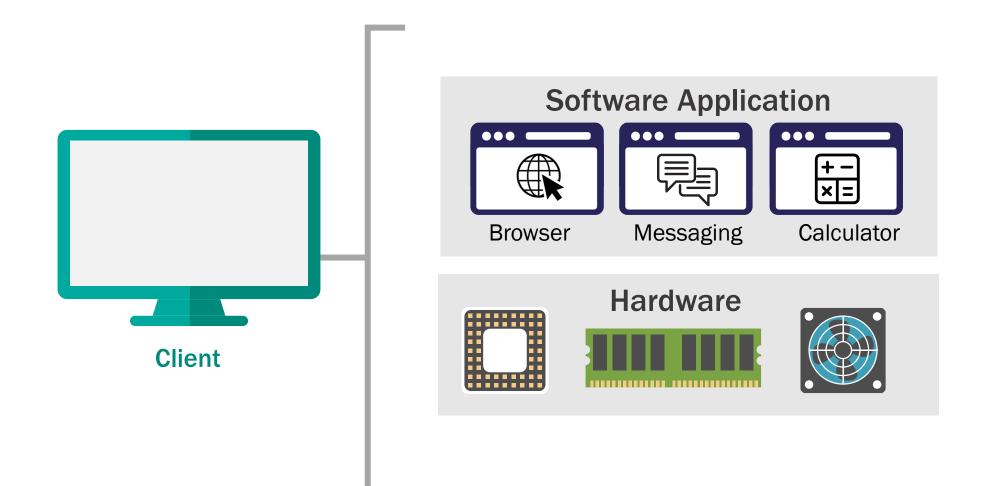


That work is not unduly burdensome.

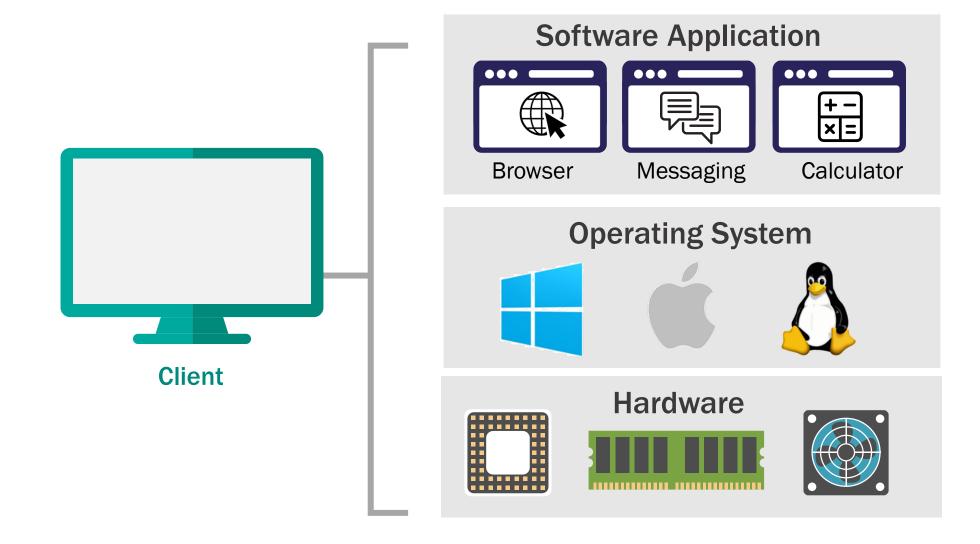


That work can be done in a way that preserves security.

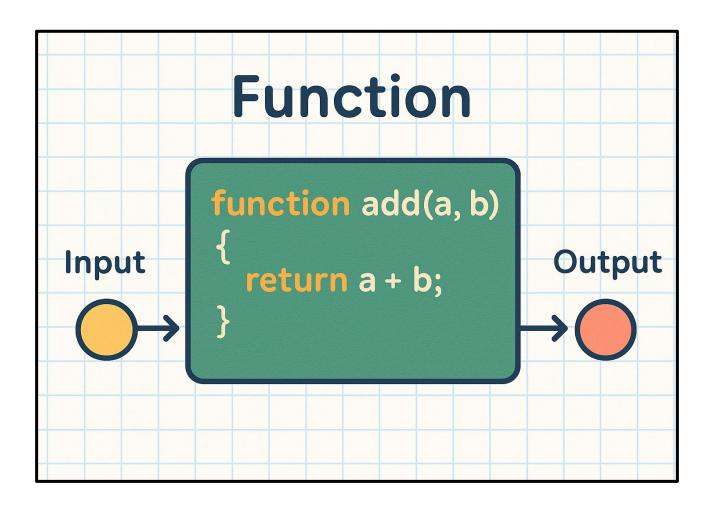
Computers: An Overview



Computers: An Overview



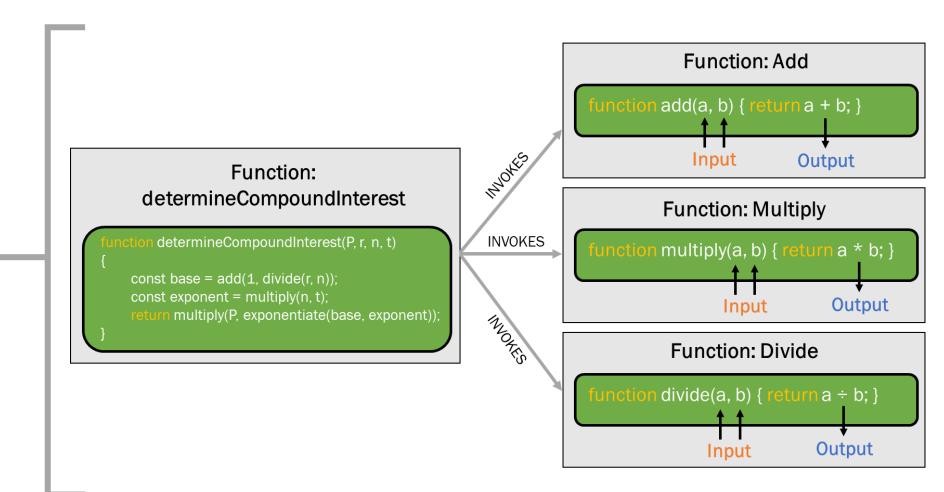
Functions Are Blocks of Code



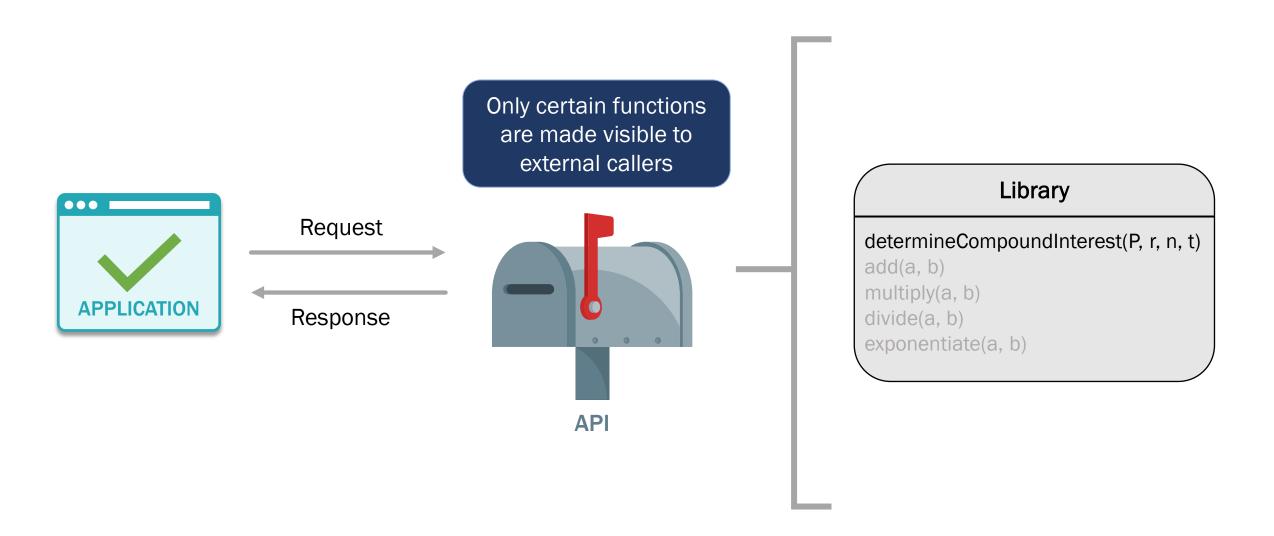
Libraries Are Collections of Functions

Library

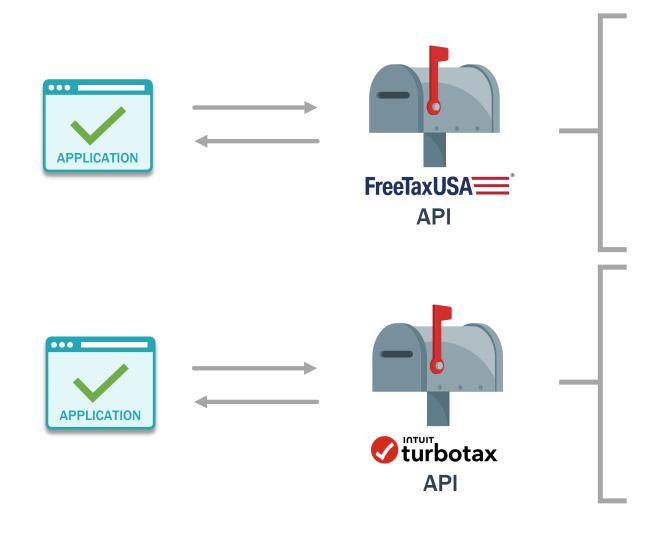
determineCompoundInterest()
add()
multiply()
divide()
exponentiate()



Application Programming Interfaces (APIs) Allow Software Programs to Talk to Each Other



Tax Example: API

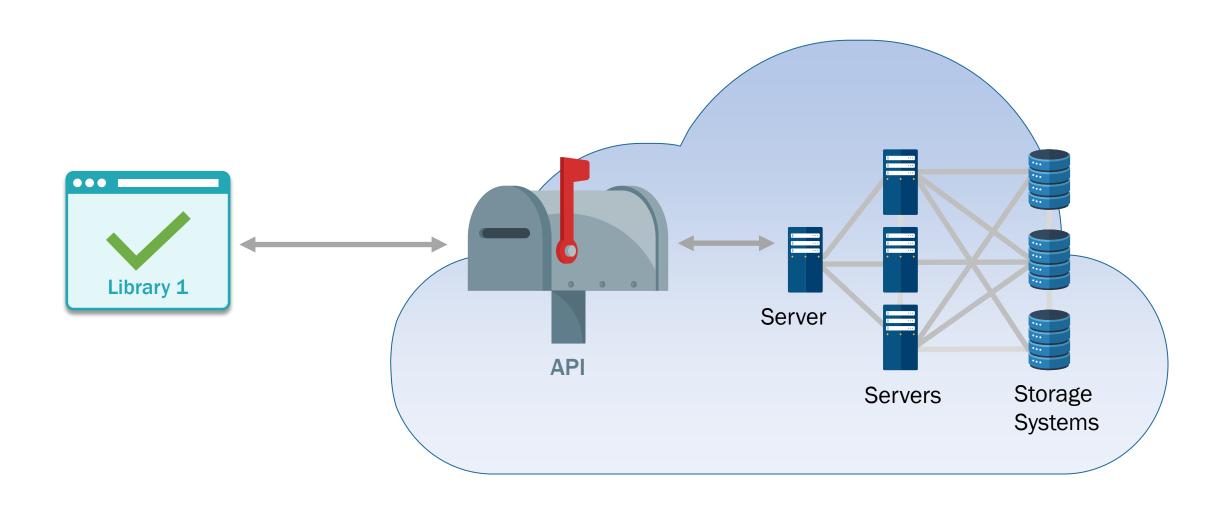


Library

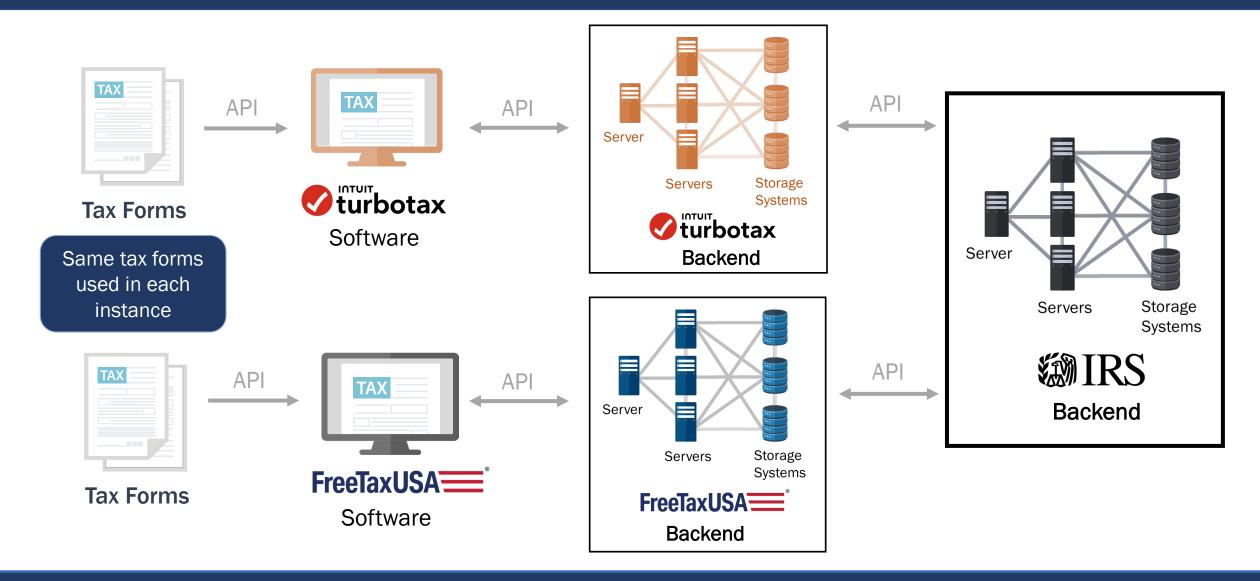
fileTaxes(salaryAndInvestmentInfo, stateOfResidence)

functionA(...)
functionB(...)
functionC(...)

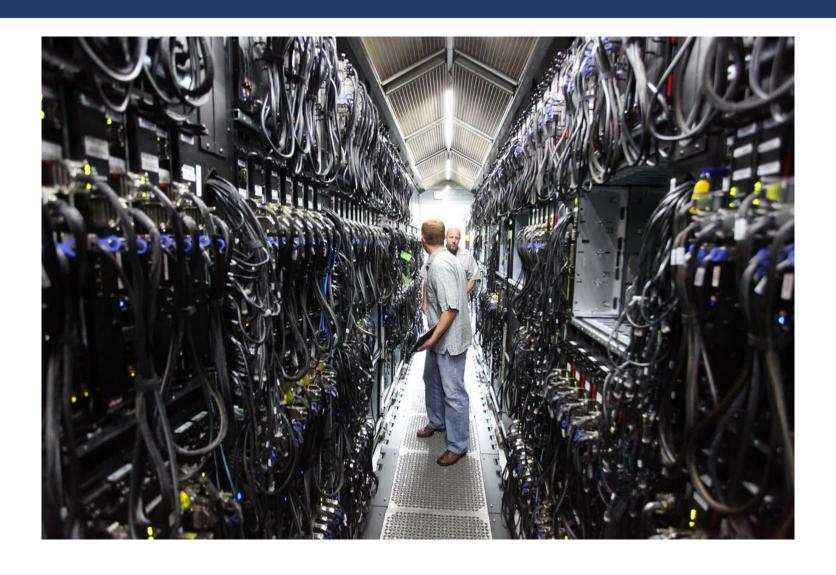
Cloud APIs Connect to Code Running in the Cloud



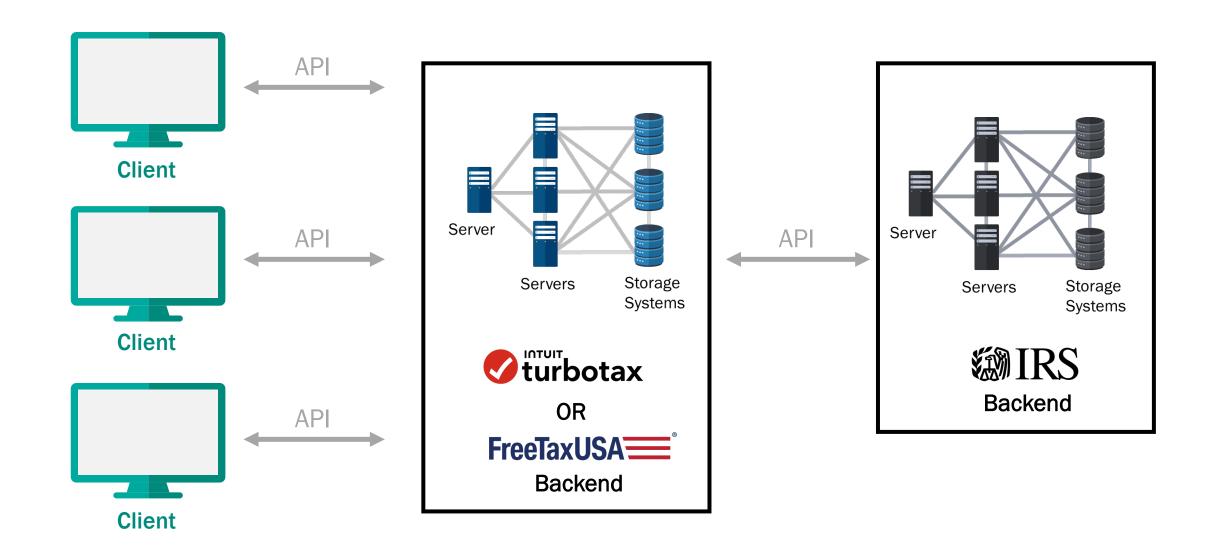
Tax Example: Distributed Systems



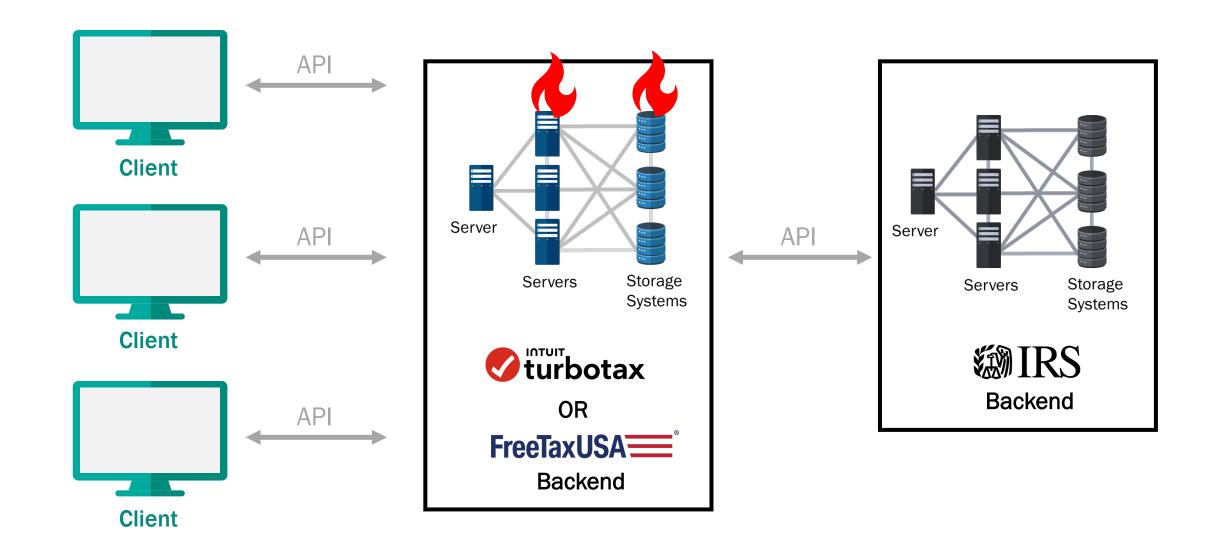
Hardware Infrastructure



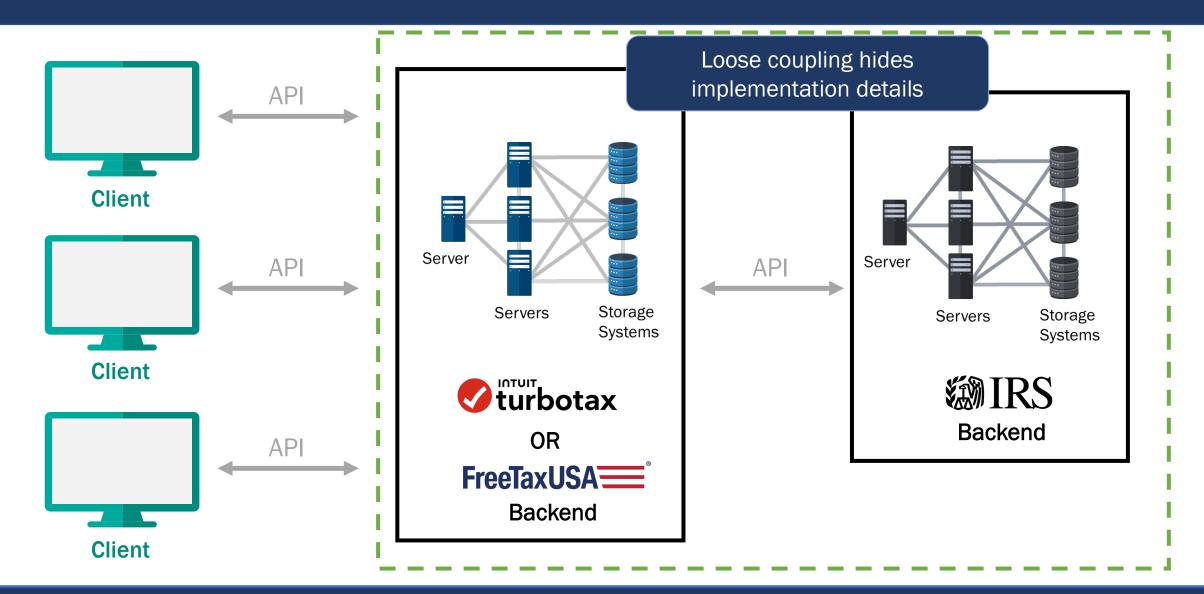
Benefits of Distributed Systems: Scalability and Efficiency



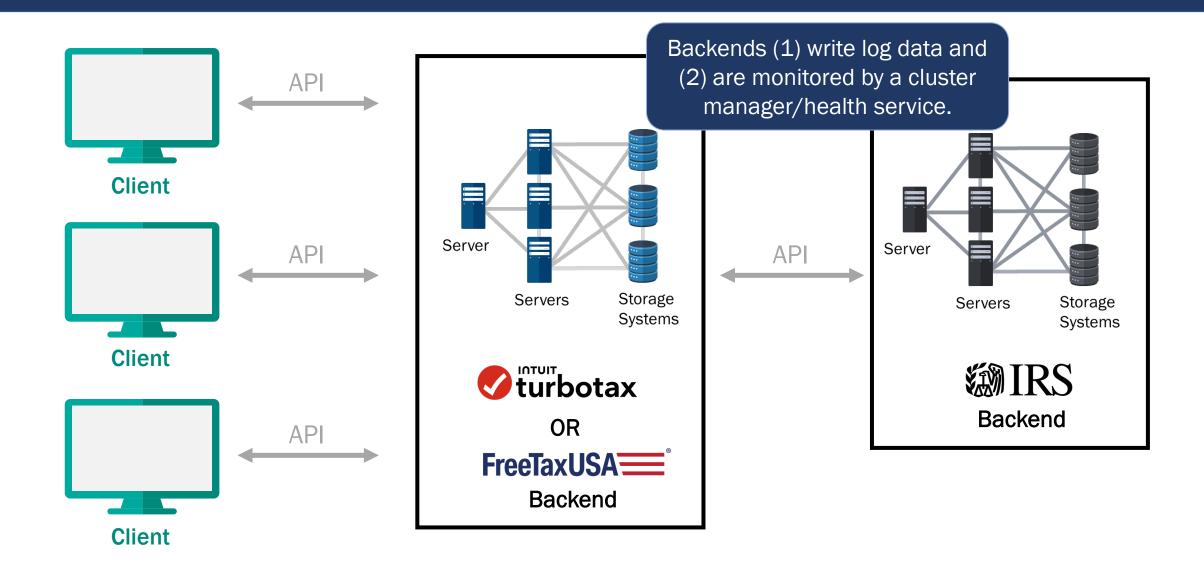
Benefits of Distributed Systems: Reliability



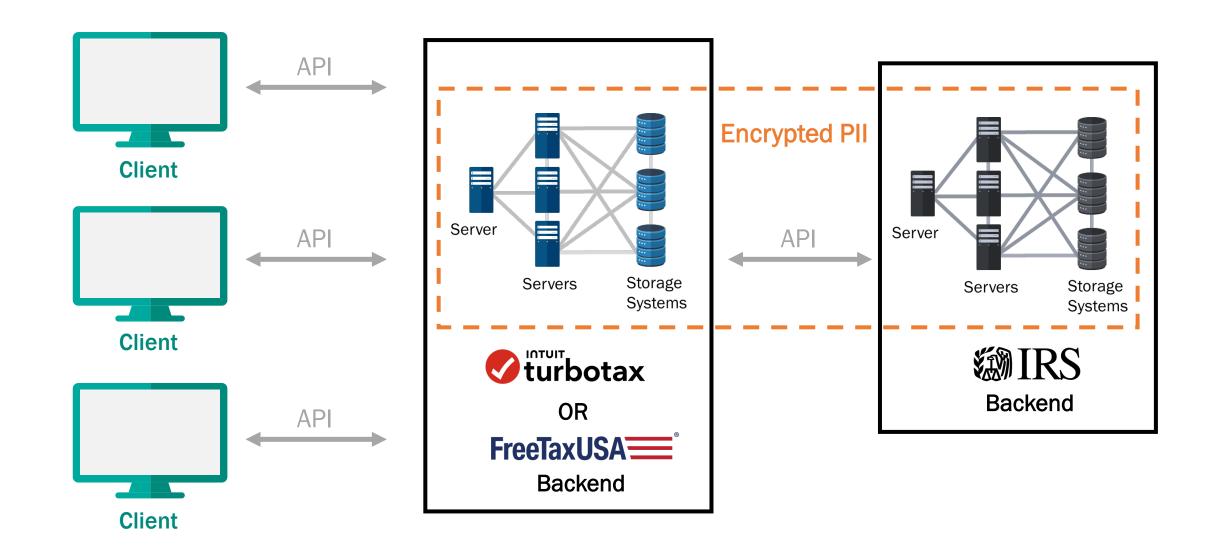
Best Practices with Distributed Systems: Loose Coupling



Best Practices with Distributed Systems: Observability



Best Practices with Distributed Systems: Security

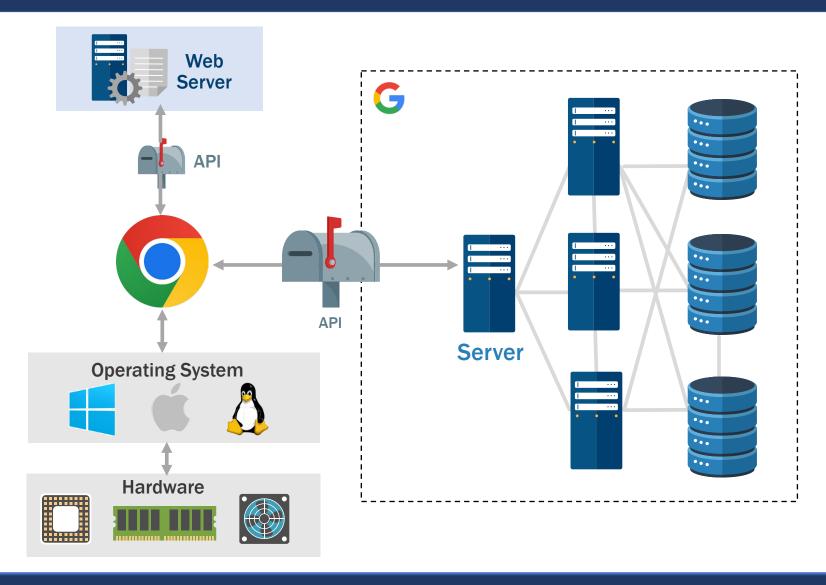


Opinion # 1: Chrome Divestiture is Feasible

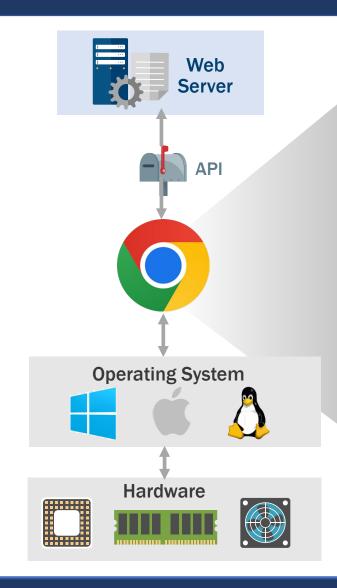
The Chrome browser is not deeply integrated with an underlying operating system.

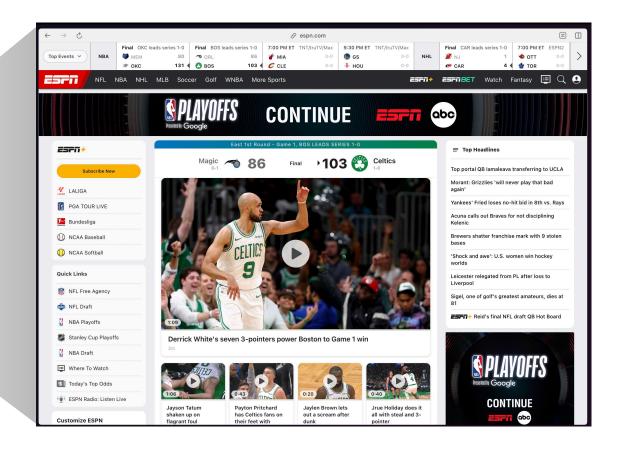
The Chrome browser is not deeply integrated with Google's back-end services.

Chrome is A Piece of Client Software

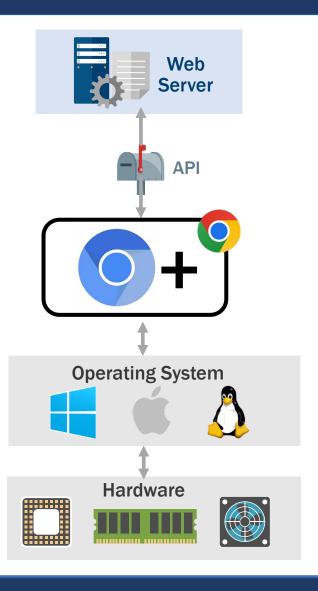


A Browser's Role

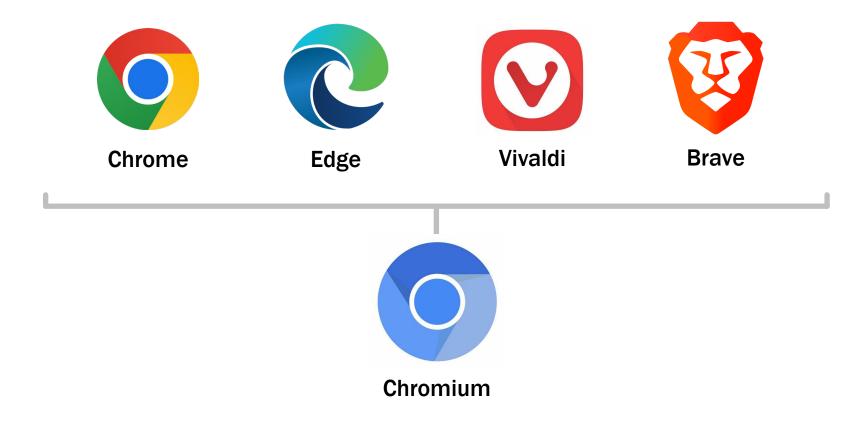




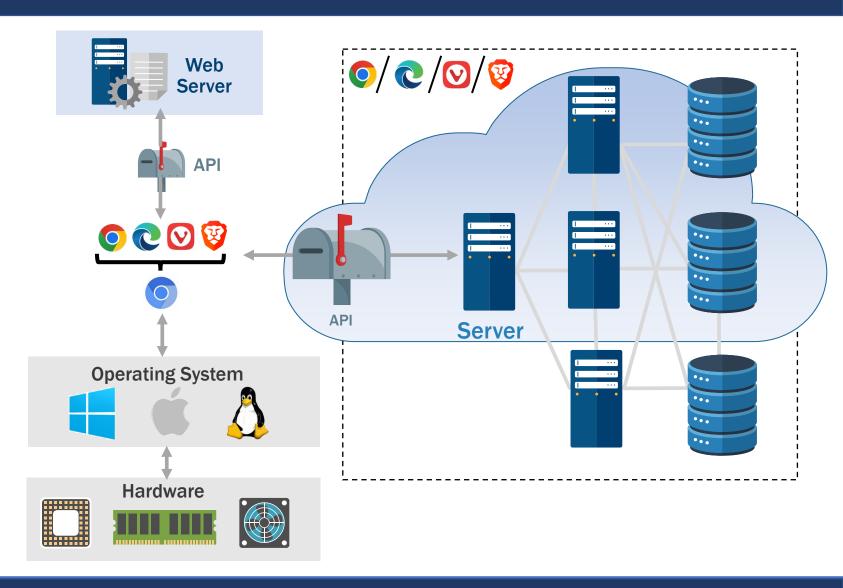
Chrome Is Built on Top of Open-Source Chromium



Other Browsers Are Built on Top of Chromium



Each Chromium-Based Browser Makes Unique Design Decisions



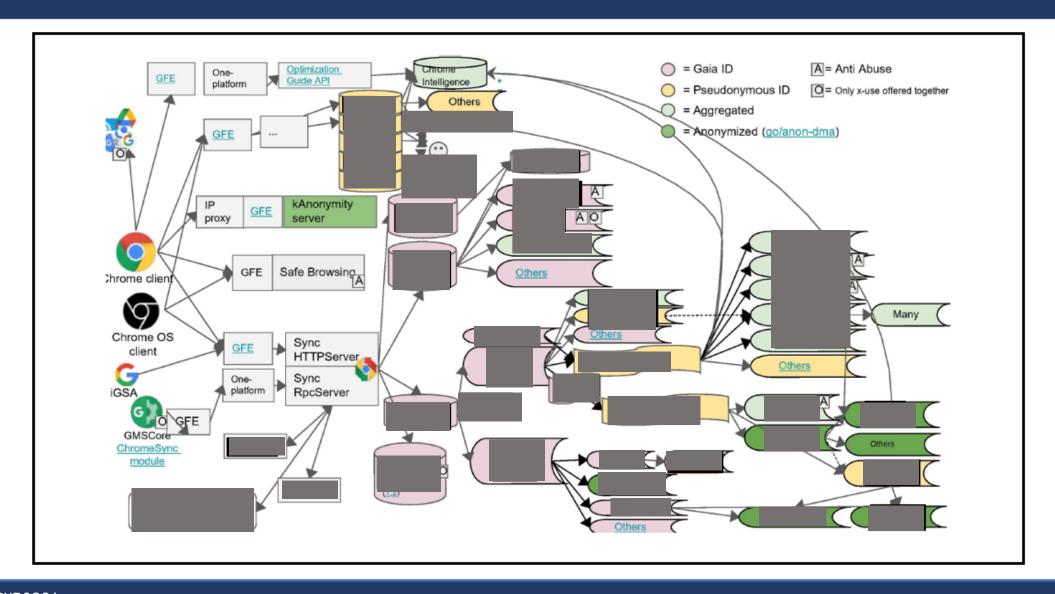
Examples of Google's Backend Services

GAIA API: Authentication service that allows users to log into Chrome.

Chrome Sync: Synchronizes Chrome data across multiple signed-in devices.

Safe Browsing API: Analyzes and identifies harmful web pages for the Safe Browsing service.

Chrome and Google's Backend



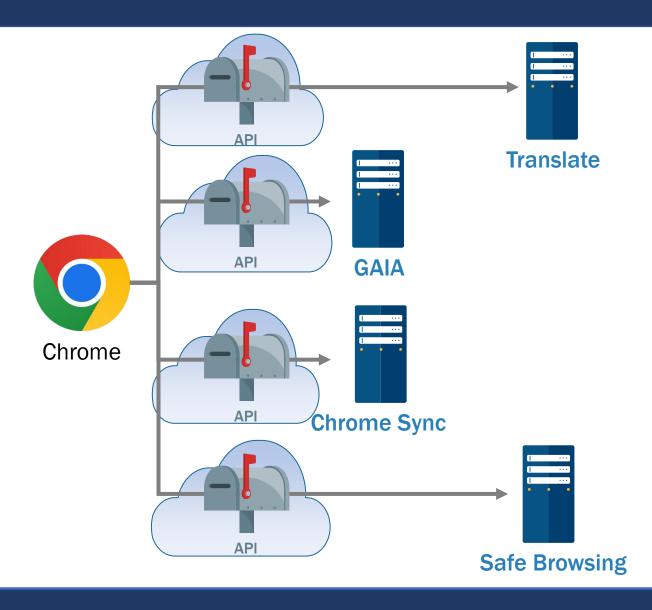
Source: PXR0321

Chrome Doesn't Need to Know Backend Implementation

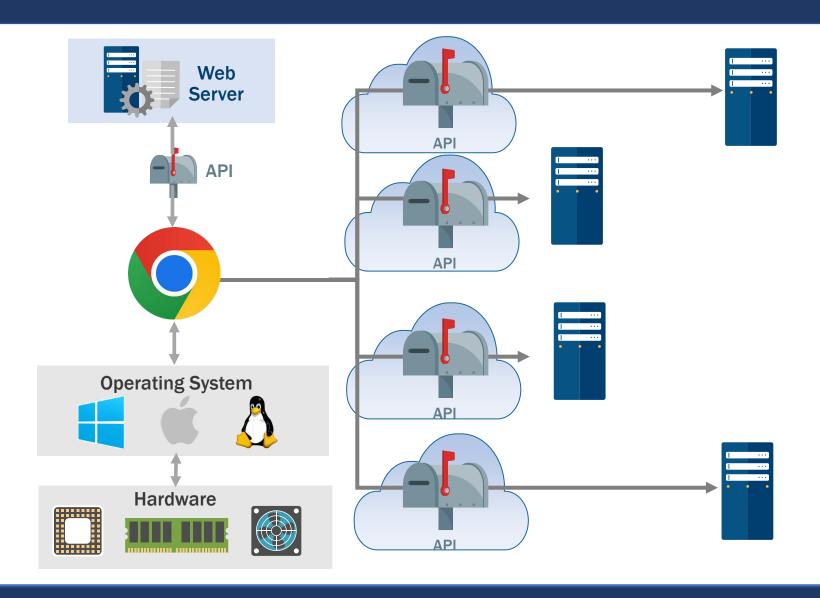


Source: PXR0321

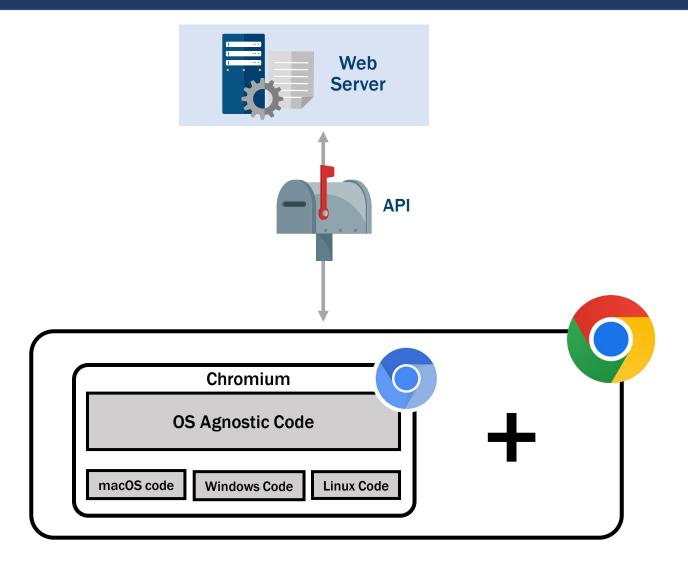
Chrome Communicates with Google Services Through Well-Defined Interfaces



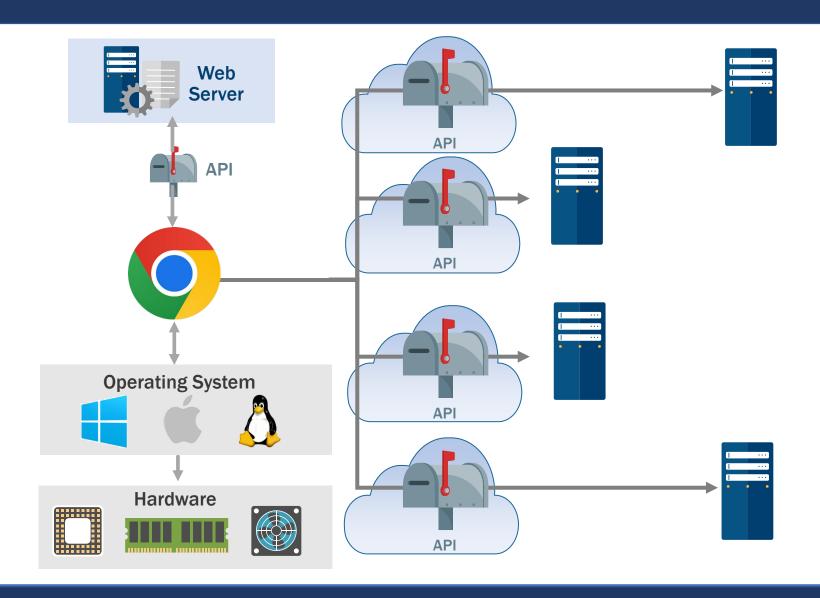
Divesting Chrome is Feasible



Divesting Chrome Would Not Impact Chrome's Ability to Retrieve and Render Web Pages



Divesting Chrome Is Feasible, Even If It May Change Chrome



The Buyer's Four Options for Each API

For each API, the new owner of Chrome can choose to:

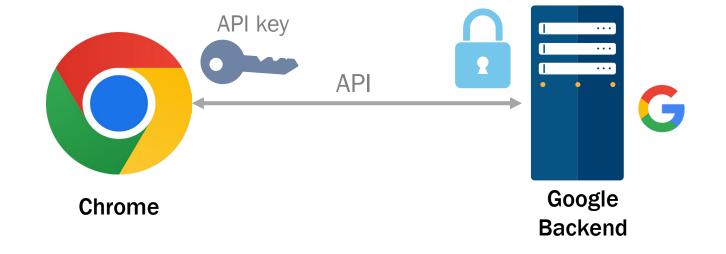
Leave unmodified the API call (e.g., purchase API keys)

Proxy the API call to a proxy server

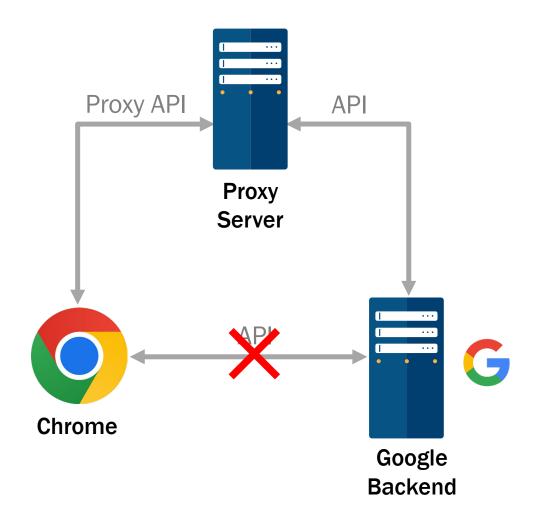
Substitute with an API call to a server run by third-party

Disable the API call

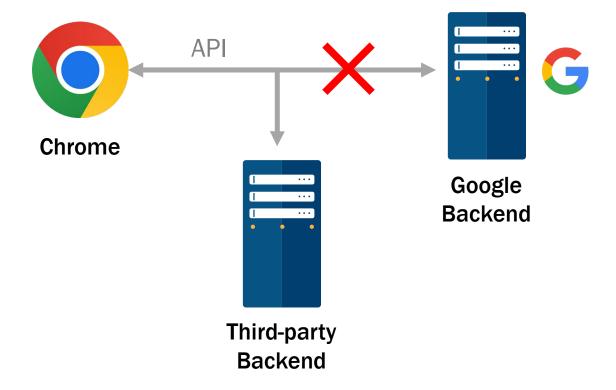
Leave Unmodified the API Call



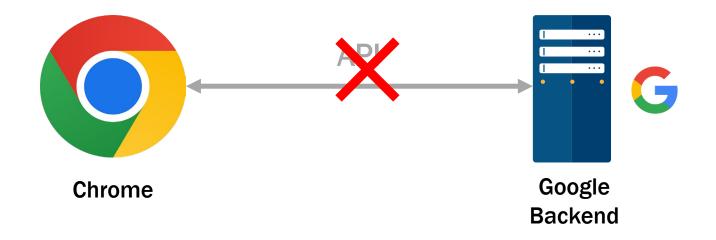
Proxy the API Call



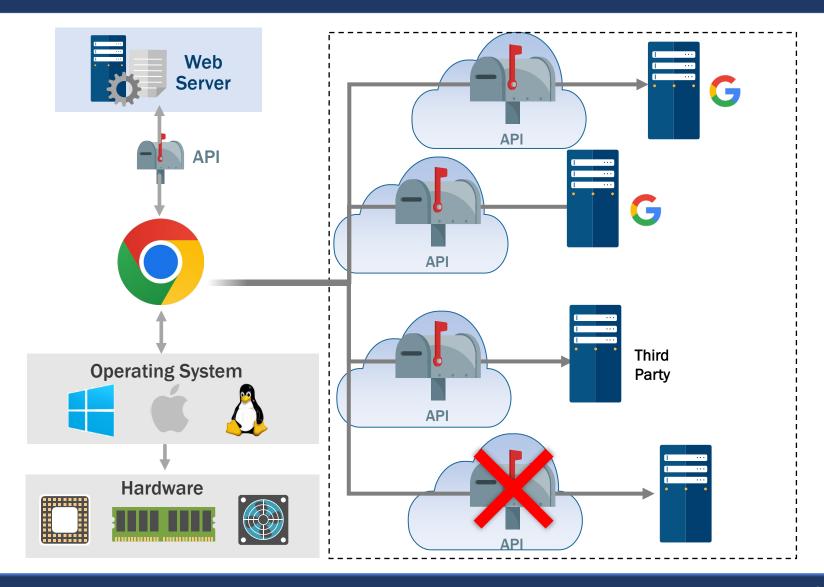
Substitute the API Call



Disable the API Call

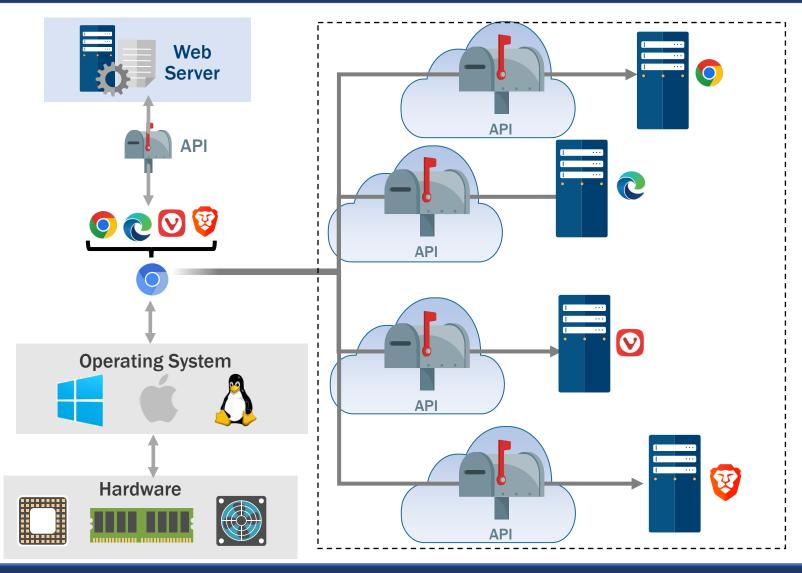


Divesting Chrome Is Feasible, Even If It Likely Changes Chrome



This Has All Happened Before





Post-Divestiture Chromium Will Receive Developer Attention









"Several leading organizations have already pledged their support for the initiative," including Google, Meta, Microsoft, and Opera. *These organizations are committed to* <u>driving innovation</u> in the Chromium ecosystem through their involvement in this initiative."

"The Supporters of Chromium-Based Browsers follow an open governance model, drawing from best practices established by other successful Linux Foundation initiatives. It prioritizes transparency, inclusivity, and community-driven development."

Source: PXR0333

This Has All Happened Before



A tool for managing the deployment of distributed systems. Kubernetes is an offshoot of Google's closed-source Borg project.



A popular closed-source browser in the 1990s. Netscape created a non-profit organization called Mozilla, who then made the code open-source, rebranded it as Firefox.



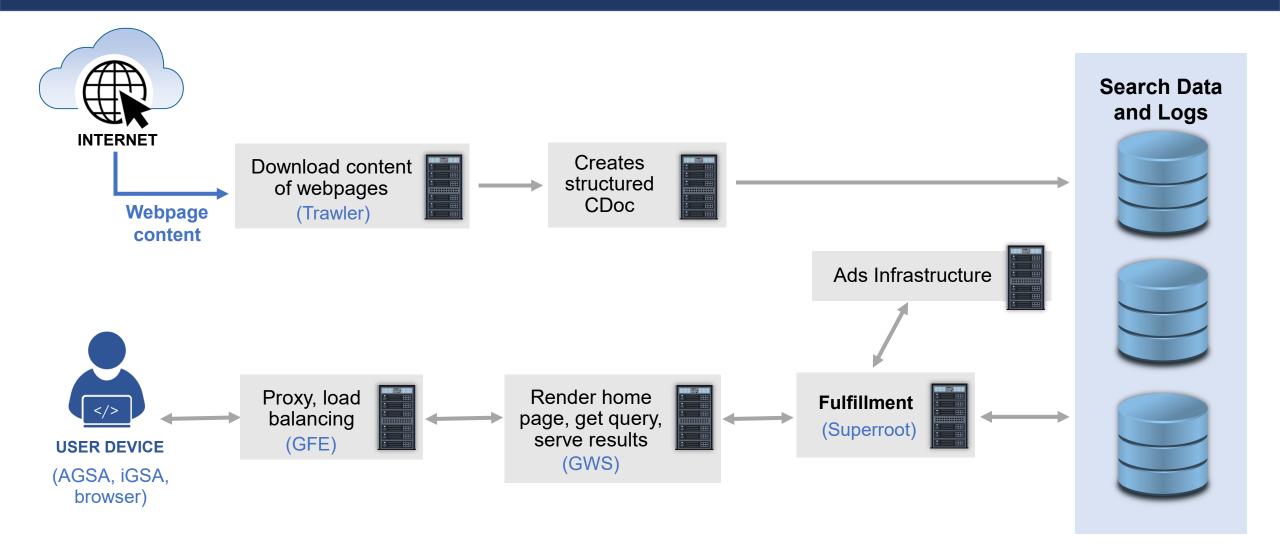
Blender was initially created as a closed-source program that was made open-source under the guidance of the non-profit Blender Foundation.

Opinion # 2: The Proposed Data Sharing and Syndication Remedies Are Feasible

Data Sharing Remedies: Google's search infrastructure is amenable to providing search engine prerequisites and ads data to Qualified Competitors.

Syndication Remedies: Google can build upon its existing search syndication systems to provide a syndicated search feed, explanatory SERP data for each syndicated query, a synthetic search feed, and a syndicated ad feed.

Overview of Google's Search Product



Source: Adapted from UPX0870 47

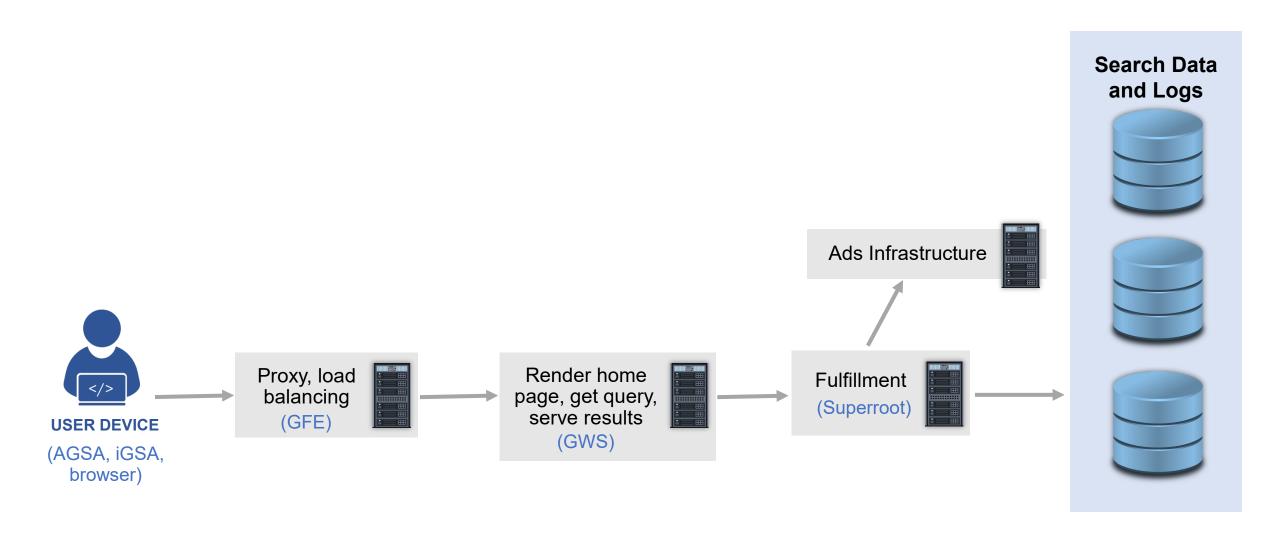
Indexing the Web





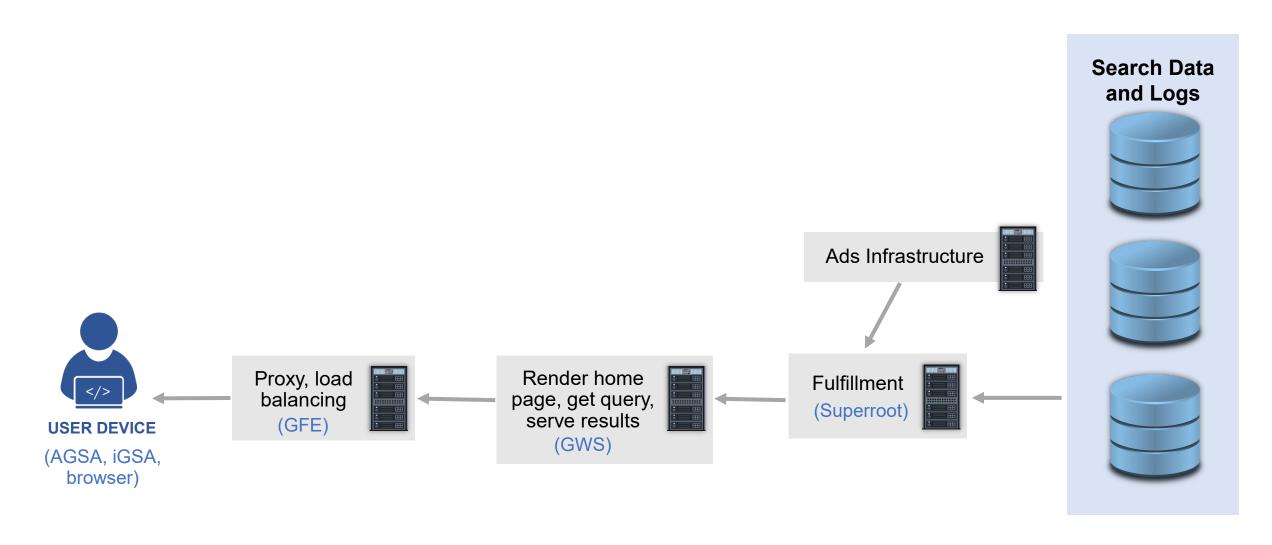
Source: Adapted from UPX0870

Part 1: Finding Results in Response to Queries



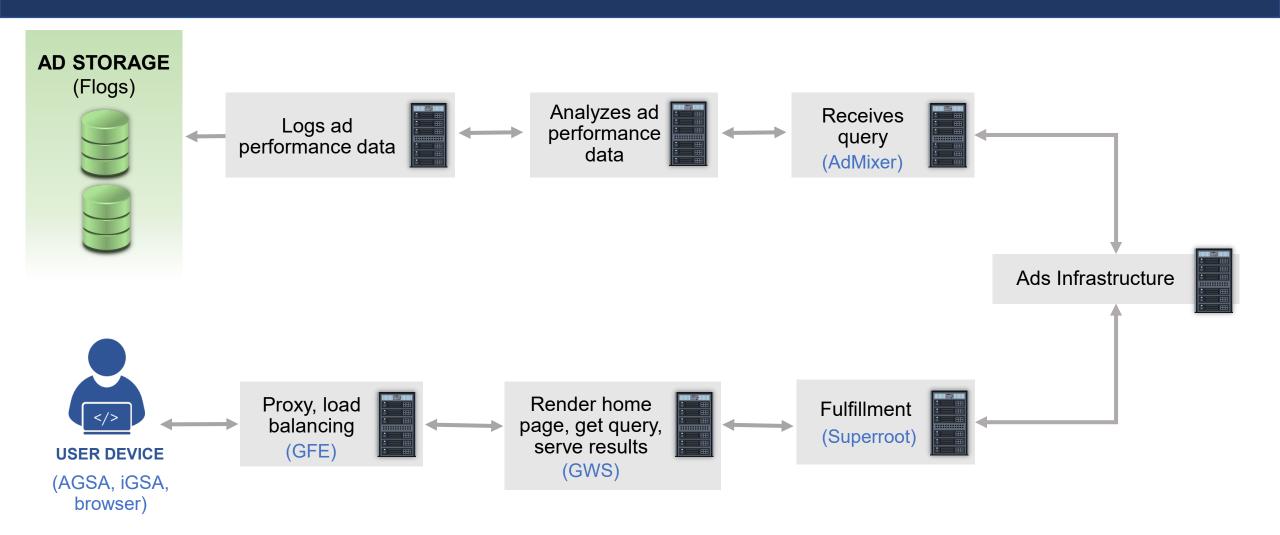
Source: Adapted from UPX0870 49

Part 2: Serving Search Results

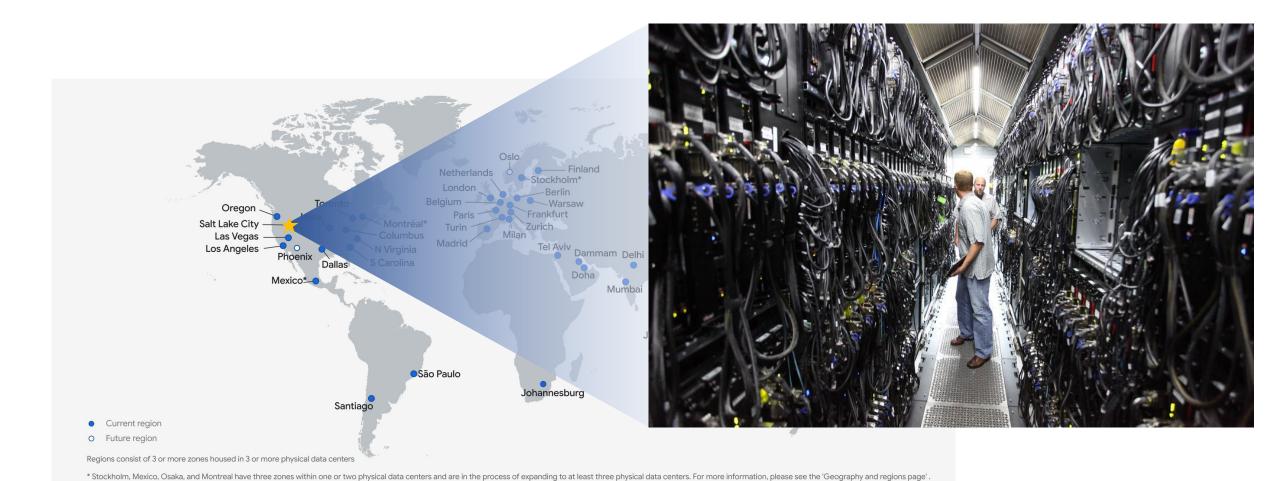


Source: Adapted from UPX0870 50

Selecting and Serving Ads

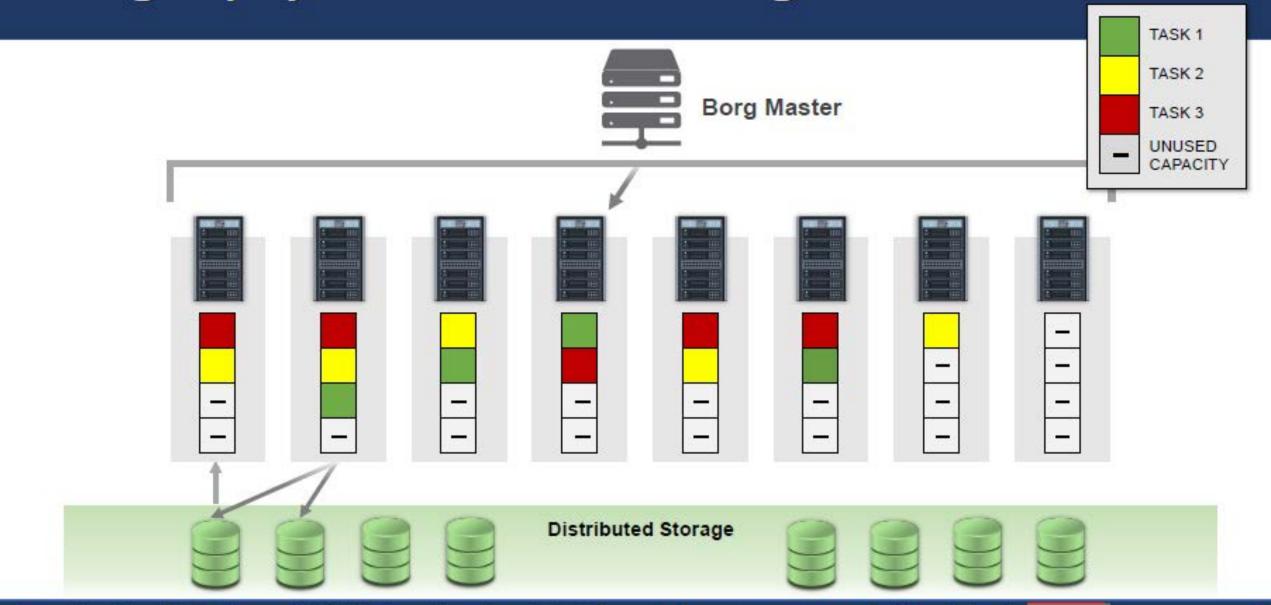


Hardware Infrastructure

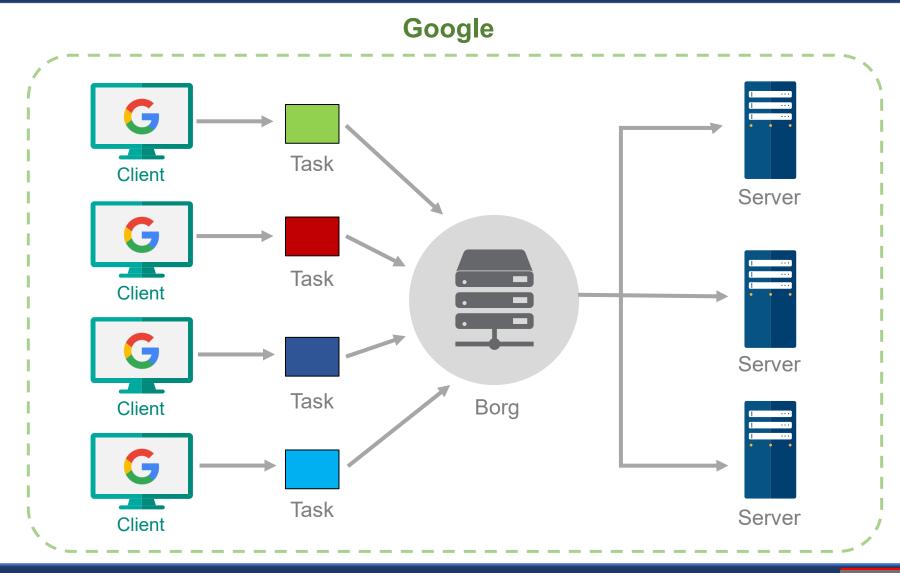


Source: https://datacenters.google/

Borg Deploys Software Across Google's Datacenters



Current Borg Example



Data Sharing Remedies

Case 1:20-cv-03010-APM Document 1184-1 Filed 03/07/25 Page 1 of 50

IN THE UNITED STATES DISTRICT COURT FOR THE DISTRICT OF COLUMBIA

UNITED STATES OF AMERICA, et al.,

Plaintiffs,

Defendant.

Case No. 1:20-cv-03010-APM
HON. AMIT P. MEHTA

GOOGLE LLC,

STATE OF COLORADO, et al.,

11

GOOGLE LLC.

Defendant.

Case No. 1:20-cv-03715-APM

HON. AMIT P. MEHTA

PLAINTIFFS' REVISED PROPOSED FINAL JUDGMENT

WHEREAS, Plaintiffs United States of America, and the States and Commonwealths of Arkansas, California, Georgia, Florida, Indiana, Kentucky, Louisiana, Michigan, Missouri, Mississippi, Montana, South Carolina, Texas, and Wisconsin, by and through their respective Attorneys General ("Co-Plaintiff States"), filed their Complaint on October 20, 2020, and their Amended Complaint on January 15, 2021;

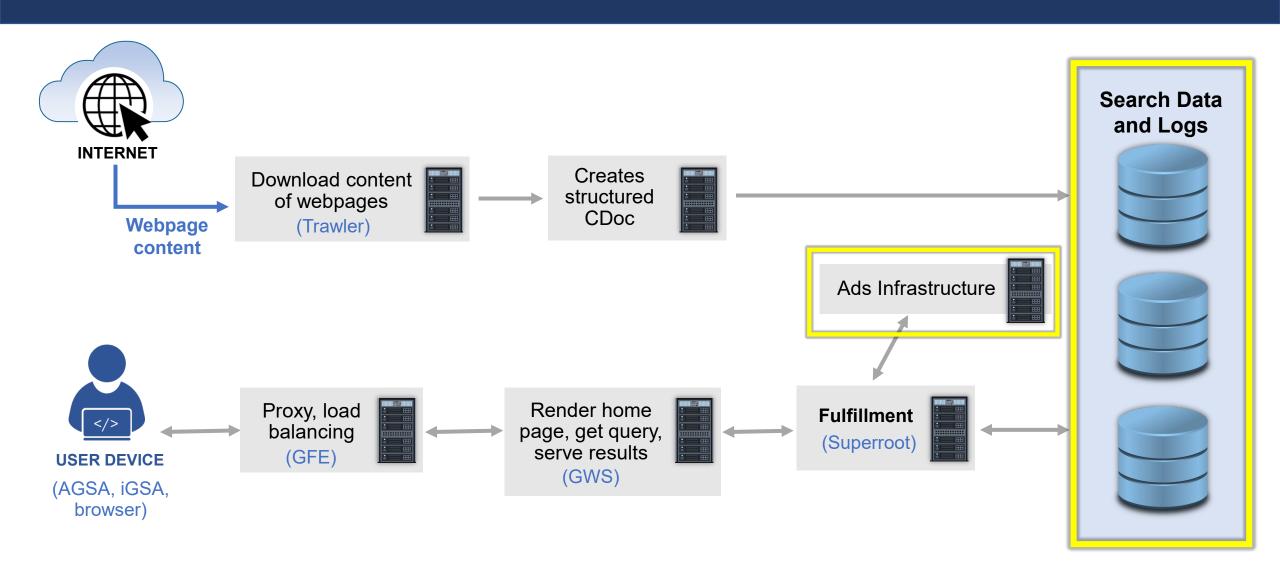
AND WHEREAS, Plaintiffs Colorado, Nebraska, Arizona, Iowa, New York, North
Carolina, Tennessee, Utah, Alaska, Connecticut, Delaware, District of Columbia, Guam, Hawaii,
Idaho, Illinois, Kansas, Maine, Maryland, Massachusetts, Minnesota, Nevada, New Hampshire,
New Jersey, New Mexico, North Dakota, Ohio, Oklahoma, Oregon, Pennsylvania, Puerto Rico,

VI. REQUIRED DISCLOSURES OF SCALE-DEPENDENT DATA NECESSARY TO COMPETE WITH GOOGLE

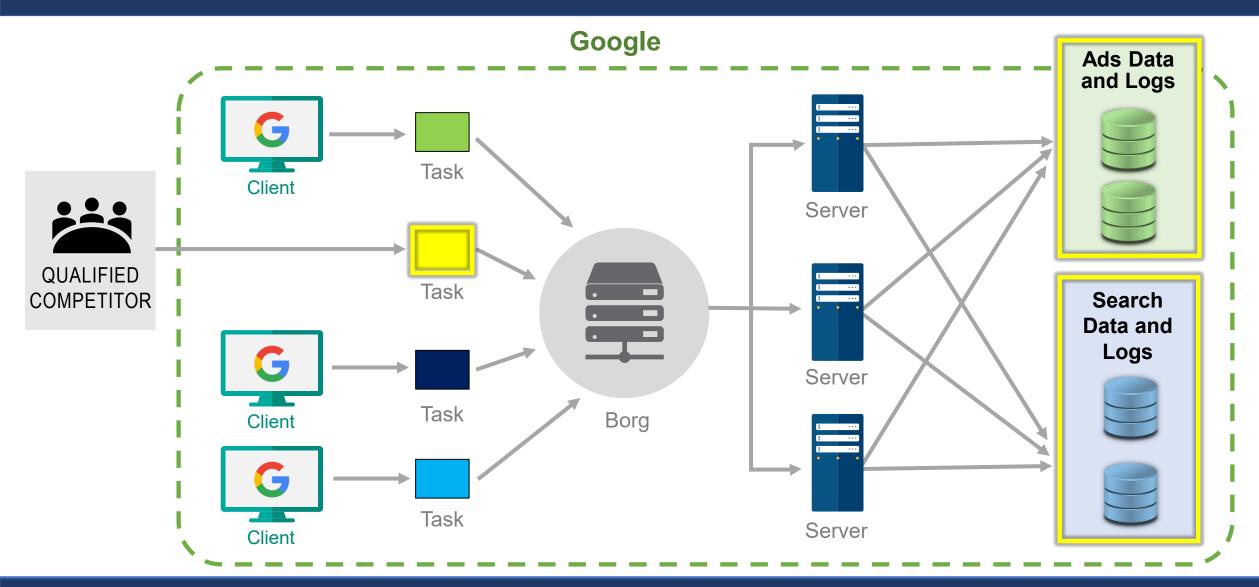
VIII. SEARCH TEXT AD TRANSPARENCY AND REDUCTION OF SWITCHING COSTS

Source: Plaintiffs' RPFJ §§ VI, VIII

Data Sharing Remedies



Google's Existing Infrastructure Supports Data Sharing Remedies



Data Syndication Remedies

Case 1:20-cv-03010-APM Document 1184-1 Filed 03/07/25 Page 1 of 50

IN THE UNITED STATES DISTRICT COURT FOR THE DISTRICT OF COLUMBIA

UNITED STATES OF AMERICA, et al.,

Plaintiffs,

GOOGLE LLC,

Defendant.

STATE OF COLORADO, et al.,

Plaintif

GOOGLE LLC,

Case No. 1:20-cv-03715-APM HON. AMIT P. MEHTA

Case No. 1:20-cv-03010-APM HON. AMIT P. MEHTA

Defendant.

PLAINTIFFS' REVISED PROPOSED FINAL JUDGMENT

WHEREAS, Plaintiffs United States of America, and the States and Commonwealths of Arkansas, California, Georgia, Florida, Indiana, Kentucky, Louisiana, Michigan, Missouri, Mississippi, Montana, South Carolina, Texas, and Wisconsin, by and through their respective Attorneys General ("Co-Plaintiff States"), filed their Complaint on October 20, 2020, and their Amended Complaint on January 15, 2021;

AND WHEREAS, Plaintiffs Colorado, Nebraska, Arizona, Iowa, New York, North

Carolina, Tennessee, Utah, Alaska, Connecticut, Delaware, District of Columbia, Guam, Hawaii,

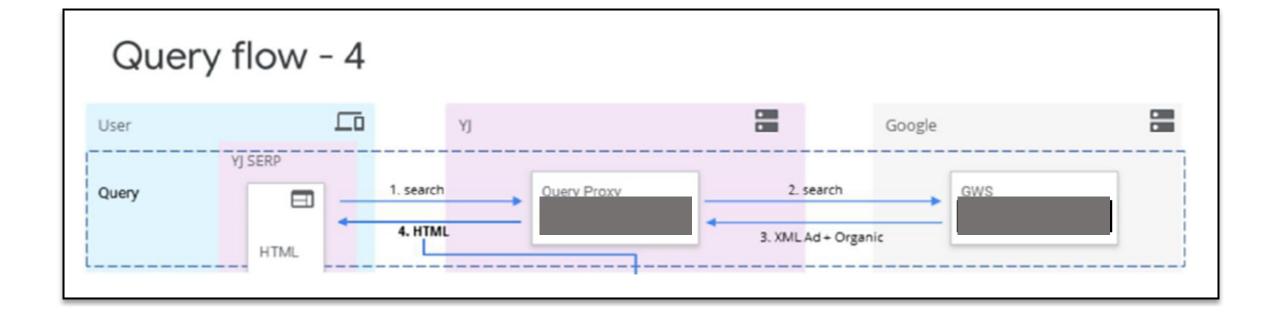
Idaho, Illinois, Kansas, Maine, Maryland, Massachusetts, Minnesota, Nevada, New Hampshire,

New Jersey, New Mexico, North Dakota, Ohio, Oklahoma, Oregon, Pennsylvania, Puerto Rico,

II. REQUIRED SYNDICATION OF SEARCH RESULTS NECESSARY TO BUILD GSE QUALITY AND SCALE OF QUALIFIED COMPETITORS

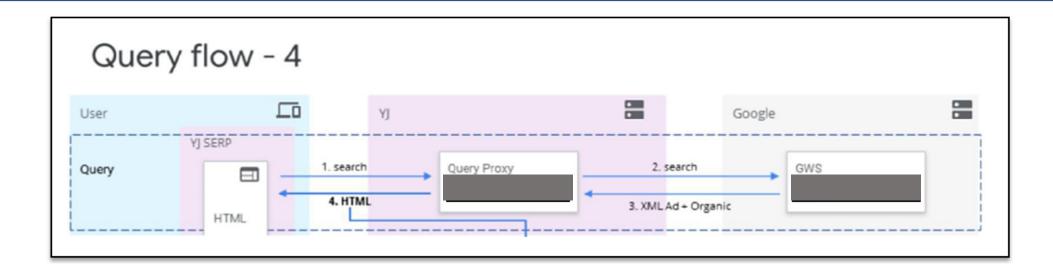
Source: Plaintiffs' RPFJ § VII 58

This Has All Happened Before



Source: PXR0318, at -083 59

There Is Technical Precedent for Syndication





Source: PXR0318, at -083 60

Professor Allan's Claims on Reverse-Engineering Are Misleading

CONTAINS HIGHLY CONFIDENTIAL INFORMATION

IN THE UNITED STATES DISTRICT COURT FOR THE DISTRICT OF COLUMBIA

United States of America et al.,

Google LLC,

Defendant.

State of Colorado et al.,

Google LLC,

Defendant.

Case No. 1:20-cv-03010-APM

Case No. 1:20-cv-03715-APM

CONTAINS HIGHLY CONFIDENTIAL INFORMATION

EXPERT REPORT OF PROFESSOR JAMES ALLAN, PH.D.

March 14, 2025

18. Given the breadth of Plaintiffs' proposals and recent advances in LLMs, implementing Plaintiffs' proposals would allow a competitor to reverse engineer essentially all of Google's search results or, at a minimum, reverse engineer Google's current and future technologies (Section IX). In particular, I find that:

Professor Allan's Claims on Reverse-Engineering Are Misleading

- 1. QCs do not want to build a perfect replica of Google.
- OpenAl and DuckDuckGo representatives have testified they do not want to replicate google.com.
- 2. A QC could not solely rely on syndicated feeds and logs to replicate a search engine.
- o **People**: engineers are needed to write code to replicate a service
- o **Time**: engineers must design, build, deploy, and test code
- o **Infrastructure**: code needs to run on physical datacenter infrastructure

Thus, QCs would be in a perpetual game of catch-up.

Professor Allan's Deposition Testimony

Q. So you're not offering an opinion on quality infrastructure latency that it would take a qualified competitor using this type of LLM or actually any – any of the reverse engineering methodologies in your report?



A. Right. All of what I'm focusing on is what is disclosed, whether directly or reverse engineering or mimicking, and what a competitor – you know, and that these are – these are complicated technologies that are being disclosed.

And I'm not saying – I'm saying the competitor could use that to improve their system. I'm not saying whether their system would then be fast, and I have not investigated how fast. Because there's too many variables in that. It's too much of a I don't know what hardware they might have. I don't know what sort of engineering expertise they have inhouse. It's just – there are too many variables there to offer that opinion.

Source: Allan Dep. at 261:24-263:2

Prohibition on Self-Preferencing in the RPFJ

Case 1:20-cv-03010-APM Document 1184-1 Filed 03/07/25 Page 1 of 50

IN THE UNITED STATES DISTRICT COURT

UNITED STATES OF AMERICA. et al..

Plaintiffs.

Case No. 1:20-cv-03010-APM

HON. AMIT P. MEHTA GOOGLE LLC,

Defendant.

STATE OF COLORADO, et al.,

GOOGLE LLC,

Defendant.

Case No. 1:20-cv-03715-APM

HON. AMIT P. MEHTA

PLAINTIFFS' REVISED PROPOSED FINAL JUDGMENT

WHEREAS, Plaintiffs United States of America, and the States and Commonwealths of Arkansas, California, Georgia, Florida, Indiana, Kentucky, Louisiana, Michigan, Missouri, Mississippi, Montana, South Carolina, Texas, and Wisconsin, by and through their respective Attorneys General ("Co-Plaintiff States"), filed their Complaint on October 20, 2020, and their Amended Complaint on January 15, 2021;

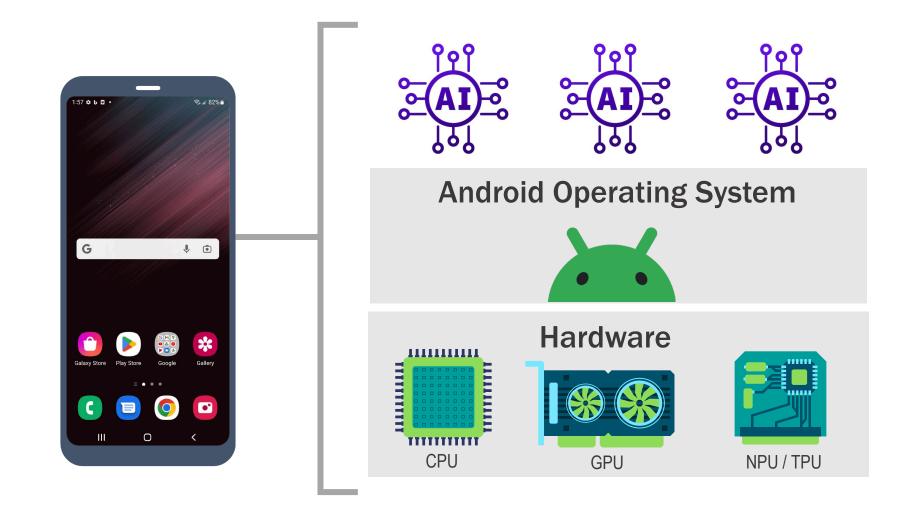
AND WHEREAS, Plaintiffs Colorado, Nebraska, Arizona, Iowa, New York, North Carolina, Tennessee, Utah, Alaska, Connecticut, Delaware, District of Columbia, Guam, Hawaii, Idaho, Illinois, Kansas, Maine, Maryland, Massachusetts, Minnesota, Nevada, New Hampshire New Jersey, New Mexico, North Dakota, Ohio, Oklahoma, Oregon, Pennsylvania, Puerto Rico

В. Self-Preferencing Prohibited: Google must not use its ownership and control of Android, or any other Google product or service, to:

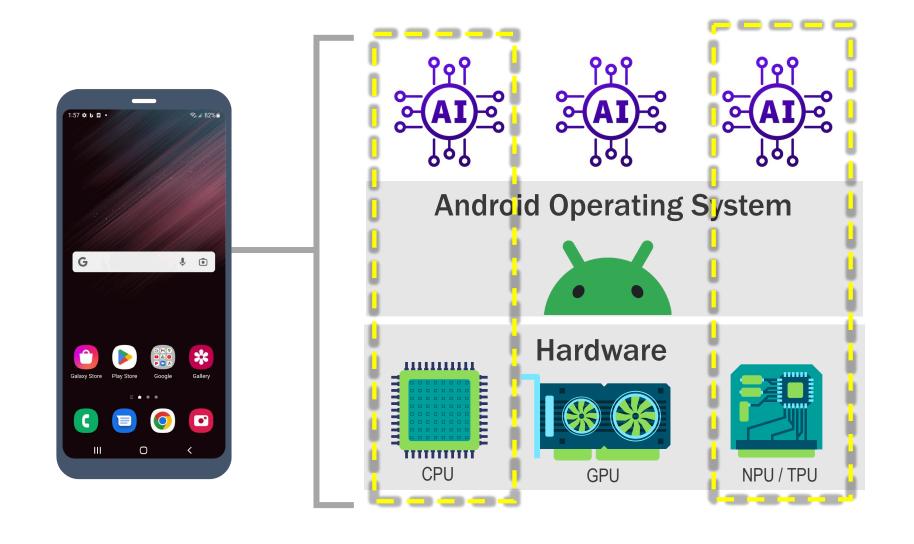
> make any GSE, Search Access Point, GenAI Product, or On-Device AI explicitly or implicitly mandatory on Android Devices, for example, by preventing interoperability between Android AICore or a Google Grounding API and Competitor products and services in the GSE or Search Text Ads markets:

Source: Plaintiffs' RPFJ & V.B 64

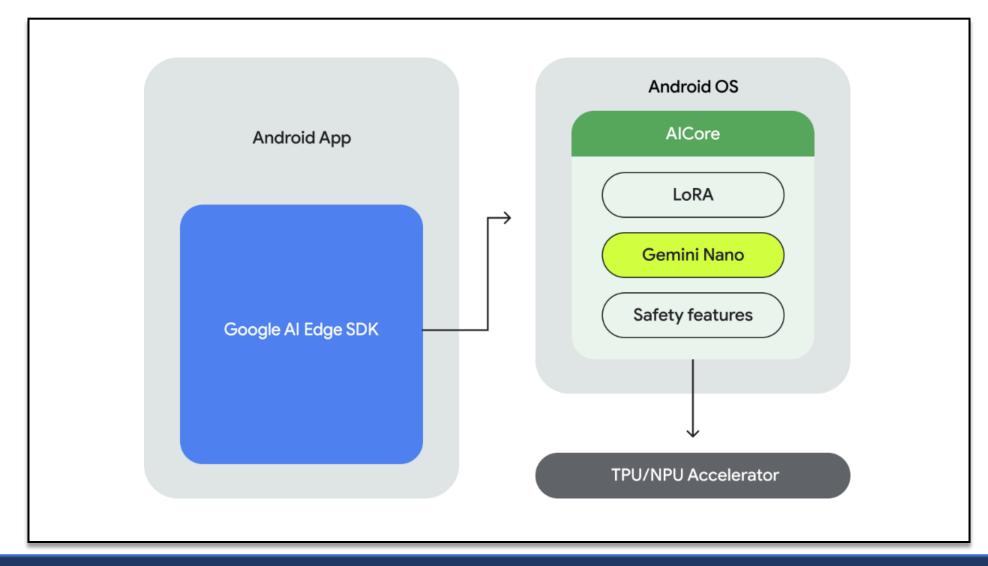
Multiple Al Models Can Run on an Android Phone



Al Models Run More Efficiently If They Run on Specialized Hardware



What Is AlCore?



Source: PXR0310 at -453

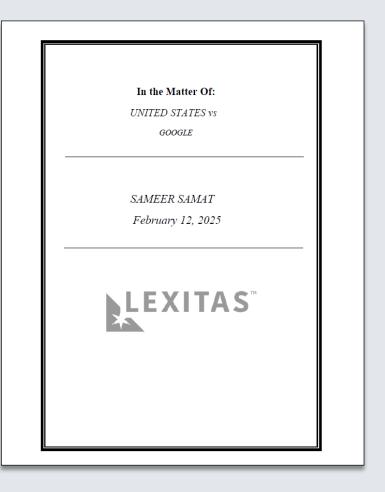
Google Is Aware of the Advantages AlCore Can Give Al Models

Google unambiguously states in public-facing technical documentation that running models on Al accelerators via AlCore is faster than running models outside of AlCore on a phone's CPUs and GPUs.

"AlCore is the new system-level capability introduced in Android 14 to provide Gemini-powered solutions for high-end devices, including integrations with the latest ML accelerators, use-case optimized LoRA adapters, and safety filters."

Source: PXR0332

Android's Lead Executive Confirms AlCore's Exclusivity

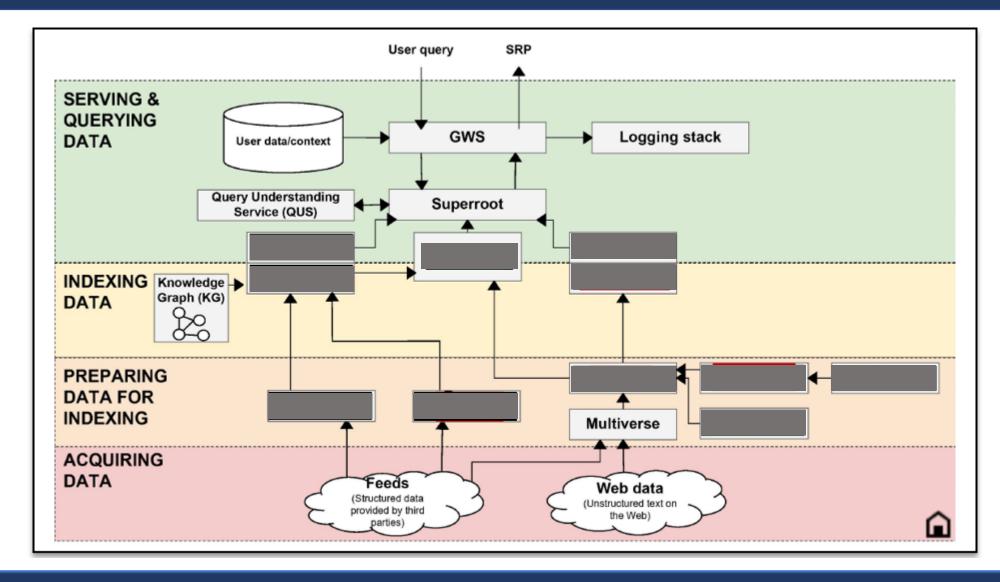


- Q. But sitting here today, are you aware of any ways for a developer to access a device resident version of Gemini Nano other than through AlCore?
- A. I'm not. But I'm also not following all the ways that they are
 they may be doing that, and there could they could exist.

* * *

- Q. Can AlCore host other device-resident LLM models?
- A. Today it cannot. We simply haven't gotten to that part of our road map yet.

Combined Search Infrastructure



Source: PXR0315 at -867

Unregulated AI Models Pose Security Risk

GUILLOTINE: Hypervisors for Isolating Malicious Als

James Mickens mickens@g.harvard.edu Harvard University

Sarah Radway sradway@g.harvard.edu Harvard University Ravi Netravali rnetravali@cs.princeton.edu Princeton University

ACM Reference Format:

James Mickens, Sarah Radway, and Ravi Netravali. 2025. GUILLO-TINIE: Hypervisors for Isolating Malicious Als. In Workshop on Hot Topics in Operating Systems (HOTOS '28), May 14–16, 2025, Banff, AB, Canada. ACM, New York, NY, USA, 9 pages. https://doi.org/10. 1145/3713082-3730931

Abstract

As AI models become more embedded in critical sectors like finance, healthcare, and the military, their inscrutable behavior poses ever-greater risks to society. To mitigate this risk, we propose Guillotine, a hypervisor architecture for sandboxing powerful AI models-models that, by accident or malice, can generate existential threats to humanity. Although Guillotine borrows some well-known virtualization techniques, Guillotine must also introduce fundamentally new isolation mechanisms to handle the unique threat model posed by existential-risk Als. For example, a rogue Al may try to introspect upon hypervisor software or the underly ing hardware substrate to enable later subversion of that control plane; thus, a Guillotine hypervisor requires careful co-design of the hypervisor software and the CPUs, RAM, NIC, and storage devices that support the hypervisor software, to thwart side channel leakage and more generally eliminate mechanisms for AI to exploit reflection-based vulnerabilities. Beyond such isolation at the software, network, and microarchitectural layers, a Guillotine hypervisor must also provide physical fail-safes more commonly associated with nuclear power plants, avionic platforms, and other types of mission-critical systems. Physical fail-safes, e.g., involving electromechanical disconnection of network cables, or the flooding of a datacenter which holds a rogue AI, provide defense in depth if software, network, and microarchitectural isolation is compromised and a rogue AI must be temporarily shut down or permanently destroyed.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for thirdparty components of this work must be honored. For all other uses, contact the owner/author(s).

HOTOS '25, May 14-16, 2025, Banff, AB, Canada © 2025 Copyright held by the owner/author(s). ACM ISBN 979-8-4007-1475-7/2025/05 https://doi.org/10.1145/2312989.2730391

1 MOTIVATION

A machine learning model tries to emulate human reasoning. To do so, a model encodes observations about training data using numerical parameters and links between those parameters. Current state-of-the-art models are so large the their internal organization is opaque to their human creators. For example, the open-source BLOOM model has 176 billion parameters [27], the open-source ILama 3.1 model has 405 billion parameters [38], and the closed-source GPT-4 models is rumored to have more than a trillion parameters [58].

Humans cannot directly understand the relationships between such vast constellations of parameters. Automated methods for understanding those relationships (and how they generate model outputs) are an active area of research. Unfortunately, such model interpretability techniques appear to be inherently fragile. Consider the task of explaining LIM inferences. The soundness of LIM interpretability techniques is vulnerable to instabilities in the underlying LIM instell [60]. For example, the fact that LIMs are sensitive to minor changes in prompt phrasing [56] can result in a model's self-reported chain-of-thought being an unfaithful representation of the model's actual reasoning process [68]. Parthermore, an LIM's tendency to hallucinate [64] can manifest itself not only in the model's answer to a question, but in the model's explanation for that answer [80].

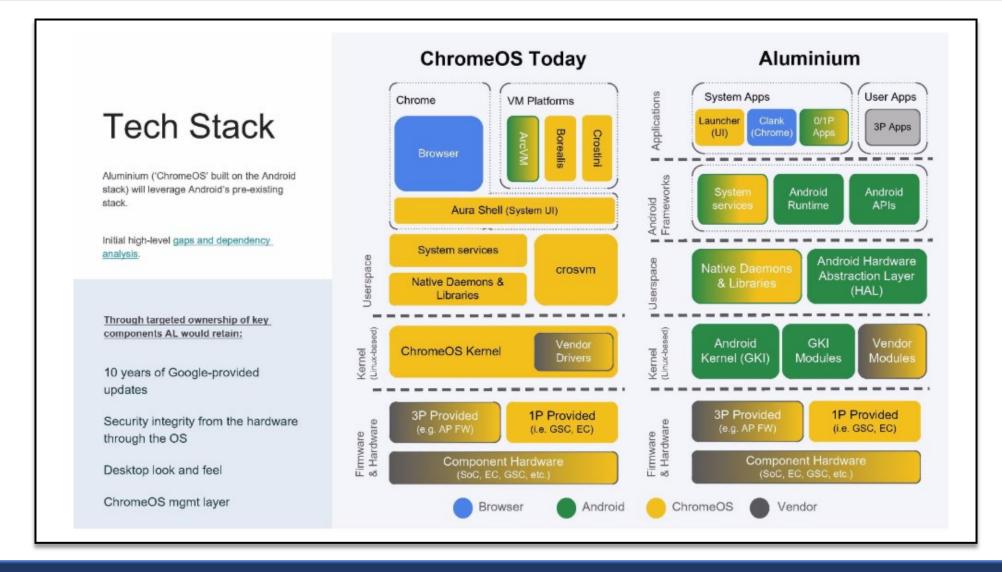
The opacity of model reasoning is troubling because models are increasingly connected to societally important infrastructure. For example, in financial settings, misbehaving models can generate huge monetary losses due to bugs. Those bugs might have been unintentionally introduced by model makers [41] or intentionally induced by adversarial examples [78]. In warlighting secancians, military leaders are already concerned that Al-governed weapons may escalate conflicts due to ignorance of geopolitical nuance [42]: these escalation problems are exacerbated when Al makes decisions to ougktky for humans to review those decisions [33].

Model alignment techniques [77] try to ensure that models adhere to human-defined behavioral norms. However (and concerningly), models can fake alignment compliance during training to later act in non-aligned ways post-deployment [21]. Researchers have also demonstrated that, in the specific context of LLMs, if model alignment does not completely climitate the possibility of an undesirable model behavior, an adversarial prompt can always elicit that behavior [77]. Thus, society faces an increasing risk that an artifical general intelligence (AGI) model which matches or exceeds human reasoning will agenerate catastrophic harms in real life.

Abstract

As AI models become more embedded in critical sectors like finance, healthcare, and the military, their inscrutable behavior poses ever-greater risks to society. To mitigate this risk, we propose Guillotine, a hypervisor architecture for sandboxing powerful AI models—models that, by accident or malice, can generate existential threats to humanity. Al-

ChromeOS and Aluminium Architecture



Source: PXR0317 at -420 73

Search Ads Machine Learning Infrastructure

[AdsML Efficiency, Convergence, and Privacy]

Authors: carlson, sunitav Last Update: October 12, 2020

See also: Trix

Description

AdsML provides components that enable highly optimized training and serving for the lat problems in the world. Our solutions play a key role in >80% of Google revenue (e.g., by models across ads revenue) and account for a large fraction of overall ML resource usage ~25% of TPU computation is through AdsML, > 2 EB of ML training data are stored in W

Key components inclu

- Flogs: High-throughput ETL system handling online feature extraction from hund sources
- Woodshed: Global, column-oriented storage system designed for ML, solving predictions, sharing features across many models, productionization, etc
 AdBrain: Training system for large-scale, TF-based models. Enables online train
- trained models like pCTR, pCVR, Ad Spam, and YouTube Viral.

 Steelmill: Serving system for ML inference in Ads, focused on high-availability,
- Steelmill: Serving system for ML inference in Ads, focused on high-availability, serving.
- <u>LegoML</u>: End-to-end platform with training focused on batch ML use cases in Ac Unsupervised learning for AdSpam, Co-trained User embeddings across Shopp

Woodshed has <u>Stacks Support Level 3</u>. Flogs, AdBrain, and Steelmill are covered by a <u>3 entry</u>. LegoML is not currently covered in Stacks, but it is believed that it already mee

There are three primary ways through which value is brought to Google by AdsML

- i. Enable *all* revenue-critical ML applications in Ads
- Contribute advanced ML infra technology/libraries back to Google (e.g., Datasync
 Leverage select AdsML components to support similar ML needs across other PA
 Woodshed for Android Security)

3 Year Vision

Our 3-year vision is organized around three strategic directions that are well aligned with goals: Efficiency, Convergence and Privacy.

Efficiency

Ads has experienced rapid growth in ML usage as we have introduced new productsfee the number of models in AdBrain has grown 8X since April 2018. We expect to see cort in the ML space because of much needed ML innovation (e.g., language and image more formats and products). Because of this, there is a large opportunity for a strategic focus cost curve for ML applications. This will allow us to confinue to enable impact with a su

Key components include:

- Flogs: High-throughput ETL system handling online feature extraction from hundreds of raw data sources
- Woodshed: Global, column-oriented storage system designed for ML, solving problems around efficiency, sharing features across many models, productionization, etc
- AdBrain: Training system for large-scale, TF-based models. Enables online training for continuously trained models like pCTR, pCVR, Ad Spam, and YouTube Viral.
- Steelmill: Serving system for ML inference in Ads, focused on high-availability, high-performance serving.
- <u>LegoML</u>: End-to-end platform with training focused on batch ML use cases in Ads. Examples:
 Unsupervised learning for AdSpam, Co-trained User embeddings across Shopping & Search models.

CONFIDENTIAL GOOG-DOJ-32659569

Source: PXR0316 at -569