# Optimizing Colocation Infrastructure for AR/VR Workloads: A Granular Bottoms - Up Forecasting Approach

## Anurag Reddy

Director Capacity Engineering, CloudFlare
Email: *anurag_reddy[at]berkeley.edu*

**Abstract:** *This paper presents a comprehensive methodology for achieving precise colocation demand forecasts through a bottoms - up approach. Colocation services play a crucial role in modern IT infrastructure, necessitating accurate predictions for effective resource allocation. The bottoms - up methodology involves a detailed analysis of individual components within the colocation landscape, such as server capacity, network capabilities, and storage requirements. By dissecting the ecosystem into fundamental elements, this approach aims to enhance the granularity of the forecasting process, providing actionable insights for organizations. The exploration delves into key elements critical to the success of bottoms - up forecasting, addressing the nuances of understanding and quantifying server capacity, network capabilities, and storage requirements. This examination contributes to a nuanced understanding of the colocation landscape, enabling informed decisions based on specific infrastructure needs. This paper unfolds a comprehensive bottoms - up methodology designed to forecast colocation demand with unparalleled precision, aligning seamlessly with the evolving needs of AR/VR landscapes. Despite the benefits, implementing bottoms - up forecasting poses challenges, including issues of data accuracy, integrating diverse components, and external factors like technological advancements and market fluctuations. The paper identifies best practices to overcome these challenges, emphasizing continuous monitoring, adaptability to evolving technologies, and flexibility in the forecasting model.*

**Keywords:** Colocation, network failover, overflow site, hard and soft constraints, interconnectivity, PoP, AR/VR applications

## 1. Introduction

Colocation demand forecasting stands as a cornerstone in the strategic capacity planning of contemporary data center management. Implementing a sophisticated bottoms - up methodology involves intricately scrutinizing forecasts at a granular level, meticulously considering individual components such as server capacity, network latency, and storage requirements, along with emerging factors like edge computing trends and artificial intelligence workloads, which collectively shape the intricate landscape of overall demand. This paper meticulously outlines the technical nuances and considerations in deploying advanced bottoms - up forecasting for colocation services, emphasizing its paramount significance in achieving unparalleled accuracy and reliability amidst the dynamic interplay of technological advancements and evolving market demands within the data center ecosystem.

## 2. Methodology

**a) Server Analysis**

1) *Historical Analysis:* Engaging in a comprehensive historical analysis of server deployments is essential to unravel intricate growth patterns, technological inclinations, and evolving client demands. By meticulously scrutinizing past data, we can unveil nuanced patterns related to the assimilation of specific server types, delve into their intricate specifications, and discern discernible trends in utilization metrics. This rigorous examination enables us to glean profound insights into the dynamic evolution of server infrastructure, facilitating data - driven decision - making that aligns with the ever - changing demands of the technological landscape

2) *Anticipated Server Types:* Incorporating technological advancements and aligning with prevailing industry trends is imperative for the development of accurate forecast models. This involves a comprehensive understanding of the evolutionary path of server technologies, encompassing enhancements in high - performance processors, augmented memory capacities, and transformations in form factors. By assimilating these insights, we aim to proactively anticipate and adapt to the forthcoming landscape of server configurations, adhering to recognized IEEE standards for precision and reliability in technological forecasting.

3) *Specifications and Workload Profiles:* A meticulous scrutiny of server specifications is executed to ascertain both power consumption and performance capabilities. Employing a bottoms - up forecasting approach mandates an exhaustive evaluation of workload profiles, aiming to estimate compute power tailored for diverse applications. This encompasses accounting for variations aligned with client requirements and dynamic shifts in industry standards. This adherence to precision in assessing server attributes aligns with recognized IEEE standards, ensuring a robust methodology in the determination of power efficiency and performance potential.

**b) Networking Infrastructure Analysis**

A robust network hardware deployment plan within a colocation facility demands an in - depth exploration of various dimensions in addition to the core aspects. Elevate the analysis with the following considerations.

1) *Bandwidth Precision:* Go beyond general assessments and conduct a granular analysis of each client's bandwidth needs. Understand the intricacies of their applications, usage patterns, and data transfer requirements. This tailored approach ensures a finely tuned network that meets the specific demands of diverse services hosted within the colocation environment.

2) *Interconnectivity Foresight:* Extend the evaluation of interconnectivity needs to encompass emerging technologies like edge computing, IoT devices, and latency sensitive applications like AR/VR. Anticipate future requirements by considering the evolving landscape of connectivity demands. This forward - looking perspective ensures the network remains agile and capable of seamlessly integrating upcoming technological advancements.

3) *Advanced Reliability Measures:* Move beyond basic reliability considerations to conduct a comprehensive risk analysis. Identify potential points of failure and proactively implement measures such as geographic diversity for redundancy, sophisticated failover mechanisms, and continuous monitoring. This approach ensures a resilient network that can withstand disruptions and offers scalability without compromising reliability.

4) *Integrated Security Protocols:* Integrate a thorough security assessment into the networking analysis. Evaluate the effectiveness of security measures, including firewalls, intrusion detection systems, and encryption protocols. Address both external threats and internal vulnerabilities to establish a robust security posture that safeguards data integrity and confidentiality.

5) *Scalability as a Core Element:* Consider scalability as a foundational element of the analysis. Anticipate the growth of client requirements and overall network demands. Develop an architecture that seamlessly scales to accommodate evolving needs without requiring substantial overhauls. This foresight is crucial for ensuring long - term sustainability and cost - effectiveness.

By incorporating these nuanced considerations, the networking infrastructure analysis becomes a comprehensive framework. This holistic approach lays the groundwork for a strategic and adaptable deployment plan, ensuring the colocation facility's network is not only reliable and secure but also future - proofed for emerging technologies and evolving client needs.
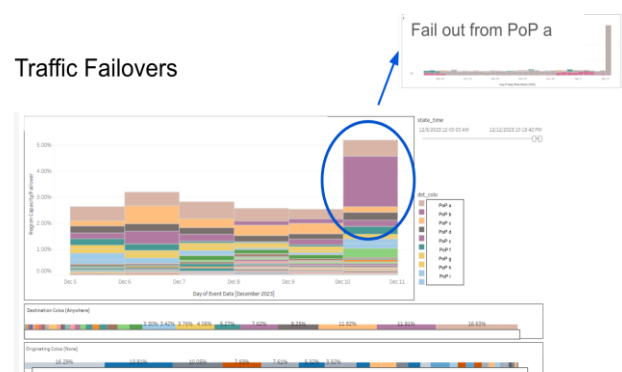
## 3. Process

### a) *Constrained Colocation Strategy*
Merging server deployments and networking infrastructure at the Point of Presence (PoP) level generates an initial "Unconstrained Colocation. " However, recognizing the need for accuracy and consumability by the sourcing team, two crucial steps are Generate a Constrained Forecast and Deliver by PoP at Net Execution PO.

1) *Within Lead - Time of Colo Expansion:* Predicting traffic failover between Points of Presence (PoPs) becomes paramount during the lead time of colocation expansion. This process entails foreseeing the necessity for traffic rerouting and implementing failover mechanisms to guarantee uninterrupted service delivery. Steps involved in this process include:
- Assessment of Expansion Timeline: Evaluate the scheduled timeline for colocation expansion. Identify key milestones and deadlines for infrastructure deployment.
- Traffic Pattern Analysis: Analyze current traffic patterns and identify potential bottlenecks. Anticipate the areas where failover mechanisms may be required during the expansion.
- Failover Mechanism Implementation: Implement failover mechanisms based on the identified areas of potential disruption. Ensure seamless transition of traffic to alternative routes.
- Regular Monitoring: Continuously monitor network performance during the lead time. Adjust failover strategies as needed based on real - time data.

2) *Outside Lead - Time:* Conducting long - term forecasting outside the lead time involves considering hard constraints that may impact colocation facility growth. This includes factors such as physical space limitations, power constraints, and regulatory restrictions. The process for forecasting traffic failover in the long term includes:
- Identification of Constraints: Identify and document physical, power, and regulatory constraints affecting colocation facilities. Understand the limitations that may impact network scalability.
- Constraint Impact Assessment: Assess how identified constraints may affect traffic failover capabilities. Quantify potential risks and vulnerabilities in the network.
- Regulatory Compliance: Ensure compliance with regulatory requirements governing colocation expansion. Plan for any necessary approvals or permits for future growth.
- Alternative Growth Strategies: Explore alternative strategies for growth considering the identified constraints. Develop contingency plans for traffic failover in case of unexpected constraints.



By following these steps, organizations can establish a comprehensive Constrained Forecast strategy tailored to the specific circumstances of datacentre colocation, whether within or outside the lead time of expansion.

### b) *Enhancing Execution through Effective Forecast Communication*

When conveying colocation requirements for a specific Point of Presence (PoP) or location, the communication should seamlessly align with the data center's topology. A mere mention of needing 100KW in X is insufficient; precision is key, specifying the exact location within all physical buildings of the Ashburn PoP where the 100KW is required. To ensure execution simplicity and strategic planning, incorporate the following essential steps:
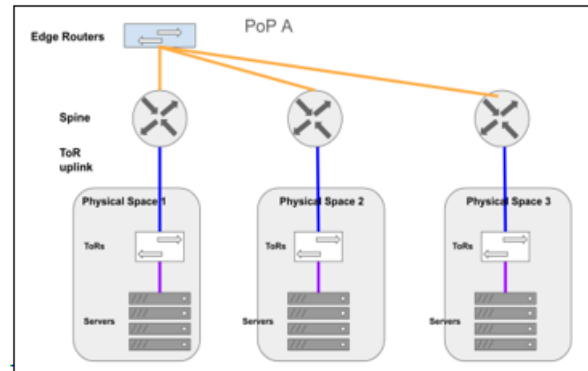
1) *Assess Current Peak Power Consumption:* Evaluate the current peak power consumption at each colocation. Identify the incremental power needs for upcoming deployments based on growth projections.
2) *Utilize Graph Breadth Traversal for Allocation:* Employ graph breadth traversal techniques to intelligently allocate the required KW within the available KW at the physical data centers. Minimize colocation costs and reduce lead time for subsequent expansions by optimizing KW distribution.
3) *Establish Fictitious Overflow Site for Excess KW:* In cases where the needed KW exceeds the total available KW, create a hypothetical overflow site strategically. Project overflow KW to provide a clear visualization of the required power at the PoP level, simplifying communication and planning.
4) *Detailed Location Specification:* Specify the exact location within all physical buildings of the PoP where the power requirement is essential. Provide comprehensive details to avoid ambiguity and facilitate precise execution.
5) *Strategic Planning with Overflow KW Projections:* Use overflow KW projections for strategic planning of the next colocation expansion.

By incorporating these essential steps, the communication of colocation needs becomes not only precise but also strategically aligned with the data center's topology. This approach ensures efficient execution, minimizes potential issues, and facilitates seamless planning for future expansions.

## 4. Conclusion

Crafting a bottoms - up forecast for colocation demand is a nuanced yet essential undertaking for optimizing resource allocation and capacity planning. Through a meticulous analysis of individual components, a deep understanding of customer needs, and a thorough consideration of scalability factors, data center managers can attain a forecast that is both precise and flexible. The adoption of best practices and the successful overcoming of challenges will empower organizations to adeptly navigate the dynamic terrain of colocation services. This document functions as a technical guide, providing valuable insights to enhance the decision - making process in forecasting colocation demand. In conclusion, the paper advocates for a forward - looking approach to colocation demand analysis. It underscores the importance of understanding current infrastructure states while anticipating future trends. Adopting the bottoms - up forecasting methodology enables organizations to strategically position themselves in the dynamic colocation market, making informed decisions aligned with evolving IT infrastructure needs.



From an application standpoint, this exploration into bottoms - up forecasting critically delves into the nuances of comprehending and quantifying server capacity, network capabilities, and storage requirements. The essence of this endeavor lies in its potential to offer a nuanced understanding of the colocation landscape tailored specifically to the intricate needs of AR/VR applications. Informed decision - making regarding infrastructure proves pivotal in influencing the performance and reliability of AR/VR experiences, positioning this methodology as an invaluable asset in the era of immersive technologies.

## References

[1] L. Wan, Y. Xu, Y. Cao, X. Cao, Y. Xia and Y. Xiong, "Research on Location and Capacity Planning of Data Centers Under New Energy Access, " 2021 IEEE 5th Conference on Energy Internet and Energy System Integration (EI2), Taiyuan, China, 2021, pp.2901 - 2906, doi: 10.1109/EI252483.2021.9713188.
[2] D. Gmach, J. Rolia, L. Cherkasova and A. Kemper, "Capacity Management and Demand Prediction for Next Generation Data Centers, " IEEE International Conference on Web Services (ICWS 2007), Salt Lake City, UT, USA, 2007, pp.43 - 50, doi: 10.1109/ICWS.2007.62.
[3] T. H. T. Le et al., "Auction Mechanism for Dynamic Bandwidth Allocation in Multi - Tenant Edge Computing, " in IEEE Transactions on Vehicular Technology, vol.69, no.12, pp.15162 - 15176, Dec.2020, doi: 10.1109/TVT.2020.3036470.
[4] Noormohammadpour, Max & Raghavendra, Cauligi. (2018). Datacenter Traffic Control: Understanding Techniques and Trade - offs. IEEE Communications Surveys & Tutorials.20.1492 - 1525.10.1109/COMST.2017.2782753.
[5] V. Tran, L. Wang and H. Chen, "A Spatial Co - location Pattern Mining Algorithm Without Distance Thresholds, " 2019 IEEE International Conference on Big Knowledge (ICBK), Beijing, China, 2019, pp.242 - 249, doi: 10.1109/ICBK.2019.00040.
[6] B. Aksanli and T. Rosing, "Providing regulation services and managing data center peak power budgets, " 2014 Design, Automation & Test in Europe Conference & Exhibition (DATE), Dresden, Germany, 2014, pp.1 - 4, doi: 10.7873/DATE.2014.156.
[7] D. Meisner and T. F. Wenisch, "Peak power modeling for data center servers with switched - mode power

supplies, " 2010 ACM/IEEE International Symposium on Low - Power Electronics and Design (ISLPED), Austin, TX, USA, 2010, pp.319 - 324, doi: 10.1145/1840845.1840911.

[8] Dongmei Huang, "DATA CENTER INFRASTRUCTURE MANAGEMENT, " in Data Center Handbook: Plan, Design, Build, and Operations of a Smart Data Center, Wiley, 2021, pp.627 - 644, doi: 10.1002/9781119597537. ch33.

[9] L. Jiadi, H. Yang, L. Huan, Z. Xinli and L. WenJing, "Research on Data Center Operation and Maintenance Management Based on Big Data, " 2020 International Conference on Computer Engineering and Application (ICCEA), Guangzhou, China, 2020, pp.124 - 127, doi: 10.1109/ICCEA50009.2020.00033.

[10] M. Wiboonrat, "Energy Management in Data Centers from Design to Operations and Maintenance, " 2020 International Conference and Utility Exhibition on Energy, Environment and Climate Change (ICUE), Pattaya, Thailand, 2020, pp.1 - 7, doi: 10.1109/ICUE49301.2020.9307075.