# NYU

# Distributed Financial Risk Assessment

Group 007

**TEAM**

Anirudh Garg - ag9563

Nikhil Kommineni - nk3853

Sai Preetham Bojja - sb9509

Sujana Maithili Chindam - sc10648

12.10.24

# **Abstract**

**Objective:** Develop a distributed framework for financial risk assessment.
**Technology:** Utilize Spark's distributed computing capabilities for efficient parallelization.
**Scale:** Process large-scale market data by simulating millions of scenarios.
**Outcome:** Generate robust and comprehensive risk metric VaR.

**All computations are run on the NYU dataproc cluster hosted on GCP.**
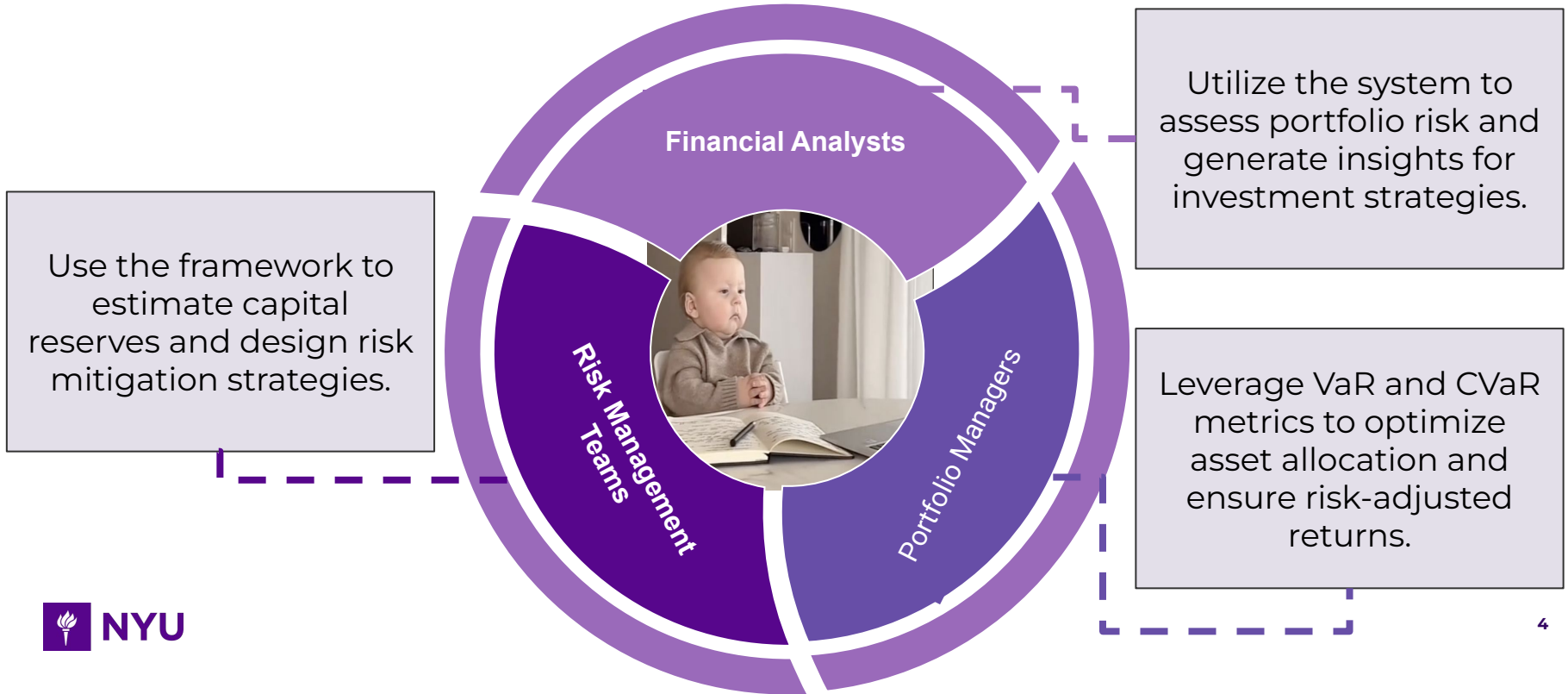
**NYU**

# Motivation

**Challenge:** Traditional risk models lack scalability and real-time adaptability.

**Need:** Accurate, scalable tools for complex financial markets.
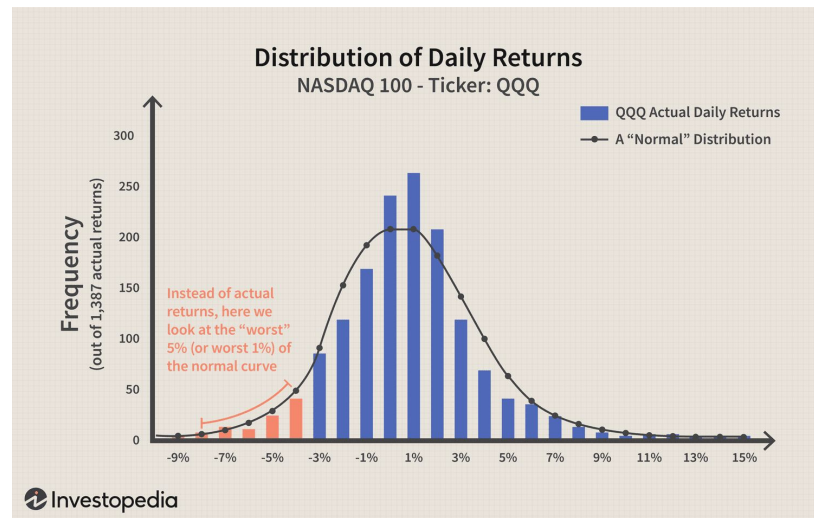
**Our work aims to:**
- Improve precision of risk analysis
- Provides actionable insights into portfolio performance and strategic decision-making
- Aid stakeholders to navigate uncertainties with confidence

# Users and Intended Beneficiaries

**Financial Analysts**

Risk Management Teams

Portfolio Managers

Utilize the system to assess portfolio risk and generate insights for investment strategies.

Use the framework to estimate capital reserves and design risk mitigation strategies.

Leverage VaR and CVaR metrics to optimize asset allocation and ensure risk-adjusted returns.

# Value at Risk (VaR)

- **A financial risk metric that estimates the maximum potential loss of an investment or portfolio over a specified time period at a given confidence level.**

- **95% VaR of $1 million means there is a 95% chance that losses will not exceed $1 million within the defined time frame.**



**Distribution of Daily Returns**
NASDAQ 100 - Ticker: QQQ

QQQ Actual Daily Returns
A "Normal" Distribution

Frequency (out of 1,387 actual returns)

Instead of actual returns, here we look at the "worst" 5% (or worst 1%) of the normal curve

Investopedia

# Conditional Value at Risk (CVaR)

- **Conditional Value at Risk (CVaR), also known as Expected Shortfall, measures the average loss in scenarios where the loss exceeds the Value at Risk (VaR) threshold.**

- **It provides a deeper insight into tail risk by focusing on the severity of extreme losses beyond the VaR, making it a crucial metric for assessing downside risk.**



**Conditional Value at Risk (CVar)**

[kən-ˈdish-nəl ˈval-(ˌ)yü ət, ˈrisk]

The amount of tail risk an investment portfolio has.

Investopedia

**NYU**

# Our Approach

| Data Collection | Data Cleaning | Analytics | Insights |
| --- | --- | --- | --- |

We collect historical financial data from a reliable source like Yahoo finance using their official python library called "yfinance".

We identify missing data points and any other inconsistencies in the collected data.

Employ different extrapolation techniques to fill in missing values in the time series data.

We run Spark jobs to compute VaR and CVaR over selected stocks using a curated list of market factors.

Run backtesting with different choices for market factors.



VAR: 0.84%



VAR: 5.4%



VAR: -4.92%

**NYU**

# DataSources

- We'll run a script for extracting data using **"wget"** command over stock symbols and google finance website.
- The stock symbols will correspond to stocks listed in US, India, Europe and Hong Kong stock markets. The list of the stocks in the countries can be extracted using the following links.
- All datasets have the same structure.

**01.NYSE**

Size: 1.1 GB - 2800 symbols - 2000 to 2024

**02.NSE**

Size: 870 MB - 2021 symbols - 2000 to 2024

**03.HKSE**

Size: 693 MB - 1800 symbols - 2000 to 2024

**04.LSE**

Size: 601MB - 1720 symbols - 2000 to 2024

NYU

# Datasets

Each dataset has two components:

a. Historical stock prices.
b. Historical Index prices.
   i. S&P 500
   ii. NASDAQ Composite
   iii. Treasury Yield 30-years (YYX)
   iv. Treasury Yield 5-years (FVX)
   v. Currency Exchange Value

## Sample Stock Data



## Sample Index Data



**NYU**

# Data Pipeline

# Results

| Stock Exchange | Value at Risk (VaR) |
|---|---|
| London Stock Exchange (LSE) | - 3.7% |
| New York Stock Exchange (NYSE) | + 2.73% |
| National Stock Exchange (NSE) | - 4.2% |
| Hong Kong Stock Exchange (HKSE) | + 2.11% |

**Indian Stock Exchange (NSE) is least risky in the short-medium term future.**

# Warren Buffett



**VaR: 4.92%**

# Bill & Melinda Gates



**VaR: 5.4%**



NYU

# George Soros



**VaR: 0.84%**

# Results

**VaR of companies listed prior to 2000: 2.7%**

**VaR of companies listed in last 4 years: 5.2%**

**It is more risky to invest in young companies compared to the established ones.**

# Challenge

## Aligning time series data with different trading calendars and handling gaps

**Standardize
Trading Calendars**

Determine the trading calendars for each dataset and establish a unified timeline based on the most restrictive trading calendar.

**Handle Non-Trading Days**

Add rows for non-trading days and fill in missing values using Forward/Backward Filling and Interpolation.

**Optimize for Large
Datasets**

Use Spark DataFrames and functions like window and lag/lead to handle missing values and align dates efficiently

**NYU**

# Challenge

## Curating the set of market factors can make or break a model

### Focus on Relevant Factors

Use statistical methods and domain insights to identify factors that significantly impact portfolio performance.

### Correlation and Multicollinearity

High correlation between factors can reduce model efficiency. Use PCA to select factors to reduce redundancy.

### Iterative Testing

Continuously refine the set of factors based on simulation results and performance metrics.

**NYU**

Distributed Financial Risk Assessment

# Challenge

## Analyzing private investment portfolio

**01** **Analyzing private investment portfolio**
- Obtain their holdings and investment details
- Obtain the respective weights of each stock in the portfolio

**02** **Handle custom weights**
- Use custom weights instead of equal-weights

The VaR assesses the potential financial risk associated with individual investors

NYU

18

# Goodness of our Results

- **Is the model internally consistent?**
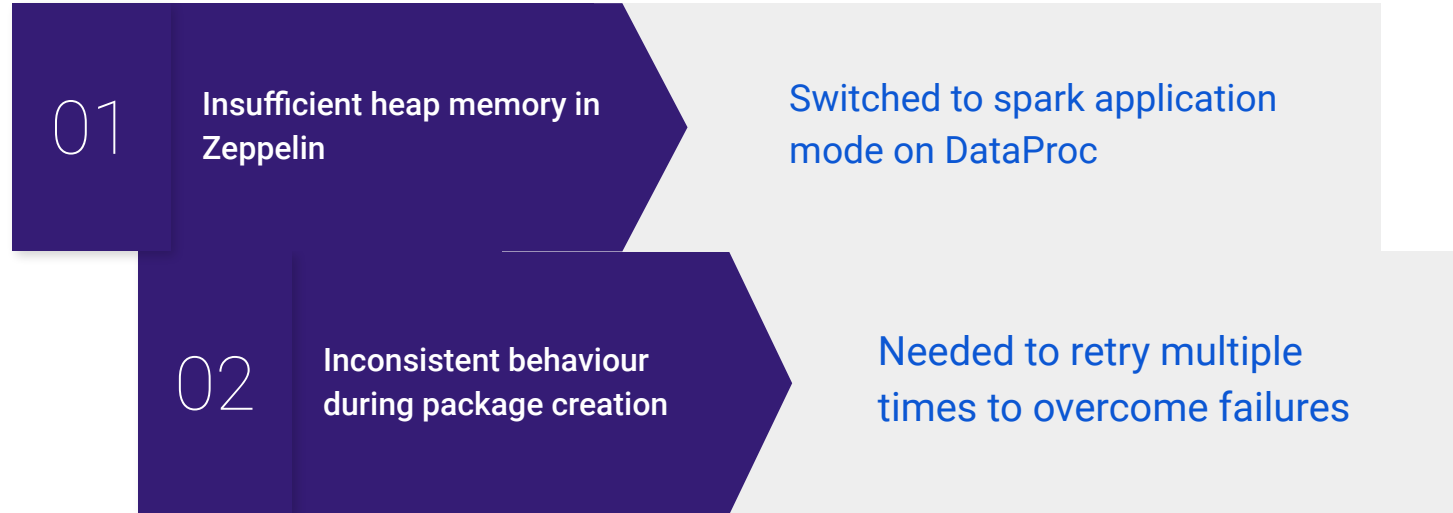  - Bootstrapping allowed us to get confidence intervals over the VaR by repeatedly sampling with replacement from the set of portfolio return results of our trials.
  - For HKSE, with high confidence we can say that VaR falls between **[0.021055, 0.021272]**.

- **How well does our model matches reality?**
  - The confidence interval does little to help us understand how well our model matches reality.
  - Kupiec's proportion-of-failures (POF) test considers how the portfolio performed at many historical time intervals and counts the number of times the losses exceeded the VaR.

**NYU**

# Obstacles

| 01 | Insufficient heap memory in Zeppelin | Switched to spark application mode on DataProc |
|---|---|---|
| 02 | Inconsistent behaviour during package creation | Needed to retry multiple times to overcome failures |

**NYU**

# ACKNOWLEDGEMENTS

- **Data Source:** We acknowledge Yahoo Finance for providing the financial data.

- **Computational Resources:** Our sincere thanks to the NYU High-Performance Computing (HPC) team for computational infrastructure

- **Academic Guidance:** We thank Prof. Yang Tang for his support throughout the course

**NYU**

# Thank You