# GeodeProposal

# Abstract

Geode is a data management platform that provides real-time, consistent access to data-intensive applications throughout widely distributed cloud architectures.

Geode pools memory (along with CPU, network and optionally local disk) across multiple processes to manage application objects and behavior. It uses dynamic replication and data partitioning techniques for high availability, improved performance, scalability, and fault tolerance. Geode is both a distributed data container and an in-memory data management system providing reliable asynchronous event notifications and guaranteed message delivery.

### Proposal

The goal of this proposal is to bring the core of Pivotal Software, Inc.'s (Pivotal) Pivotal GemFire codebase into the Apache Software Foundation (ASF) in order to build a vibrant, diverse and self-governed open source community around the technology. Pivotal will continue to market and sell Pivotal GemFire based on Geode. Geode and Pivotal GemFire will be managed separately. This proposal covers the Geode source code (mainly written in Java), Geode documentation and other materials currently available on GitHub.

While Geode is our primary choice for a name of the project, in order to facilitate PODLINGNAMESEARCH we have come up with two alternatives:

- Haptic
- FIĠ

### Background

GemFire is an extremely mature and robust product that can trace its legacy all the way back to one of the first Object Databases for Smalltalk: GemStone. The GemFire code base has been maintained by the same group of engineers as a closed source project. Because of that, even though the engineers behind GemFire are the de-facto knowledge leaders for distributed in-memory management, they have had little exposure to the open source governance process. The original company developing GemStone and GemFire was acquired by VMWare in 2010 and later spun off as part of Pivotal Software in 2013. Today GemFire is used by over 600 enterprise customers. An example deployment includes China National Railways that uses Pivotal GemFire to run railway ticketing for the entire country of China with a 10 node cluster that manages 2 gigabytes "hot data" in memory, and 10 backup nodes for high availability and elastic scale.

### Rationale

Modern-day data management architectures require a robust in-memory data grid solution to handle a variety of use cases, ranging from enterprise-wide caching to real-time transactional applications at scale. In addition, as memory size and network bandwidth growth continues to outpace those of disk, the importance of managing large pools of RAM at scale increases. It is essential to innovate at the same pace and Pivotal strongly believes that in the Big Data space, this can be optimally achieved through a vibrant, diverse, self-governed community collectively innovating around a single codebase while at the same time cross-pollinating with various other data management communities. ASF is the ideal place to meet these ambitious goals.

# Initial Goals

Our initial goals are to bring Geode into the ASF, transition internal engineering processes into the open, and foster a collaborative development model according to the "Apache Way." Pivotal plans to develop new functionality in an open, community-driven way. To get there, the existing internal build, test and release processes will be refactored to support open development.

# **Current Status**

Currently, the project code base is licensed for evaluation purposes and is available for download from Pivotal.io (https://network.pivotal.io/products/projectgeode). The documentation and wiki pages are available as public GitHub repositories under Project Geode organization on GitHub (https://github.com /project-geode). Although Pivotal GemFire was developed as a proprietary, closed-source product, the internal engineering practices adopted by the development team lend themselves well to an open, collaborative and meritocratic environment.

The Pivotal GemFire team has always focused on building a robust end user community of paying and non-paying customers. The existing documentation along with StackOverflow and other similar forums are expected to facilitate conversions between our existing users so as to transform them into an active community of Geode members, stakeholders and developers.

#### Meritocracy

Our proposed list of initial committers include the current GemFire R&D team, Pivotal Field Engineers, and several existing customers and partners. This group will form a base for the broader community we will invite to collaborate on the codebase. We intend to radically expand the initial developer and user community by running the project in accordance with the "Apache Way". Users and new contributors will be treated with respect and welcomed. By participating in the community and providing quality patches/support that move the project forward, they will earn merit. They also will be encouraged to provide non-code contributions (documentation, events, community management, etc.) and will gain merit for doing so. Those with a proven support and quality track record will be encouraged to become committers.

### Community

If Geode is accepted for incubation, the primary initial goal will be transitioning the core community towards embracing the Apache Way of project governance. We would solicit major existing contributors to become committers on the project from the start.

#### **Core Developers**

While a few core developers are skilled in working in openly governed Apache communities. Most of the core developers are currently NOT affiliated with the ASF and would require new ICLAs before committing to the project.

#### Alignment

The following existing ASF projects can be considered when reviewing Geode proposal:

Apache Hadoop is a distributed storage and processing framework for very large datasets focusing primarily on batch processing for analytic purposes. Geode is a data management platform that provides real-time, consistent, and transactional access to data-intensive applications. Our roadmap includes plans to provide close integration with HDFS.

Apache HBase offers tabular data stored in Hadoop based on the Google Bigtable model. HBase uses a key-based partitioning scheme and column family data model that can work well for scan intensive workloads but is not as broadly applicable as the rich object model, OQL querying, and hash partitioning provided by Geode. Geode will use the HFile format for storing data in HDFS.

Apache Spark is a fast engine for processing large datasets, typically from a Hadoop cluster, and performing batch, streaming, interactive, or machine learning workloads. Geode supports high throughput streaming ingest application patterns and data parallel algorithms. Geode also ensures that data is highly available through redundancy while Spark recovers from faults using RDD lineage recomputation. Our roadmap includes plans for providing integration with the Spark platform.

Apache Ignite (incubating) offers distributed in-memory processing capabilities which in some ways overlap with Geode. However, Geode has been in this field for more than 10 years and the product API's, design, and implementation details are quite different. In addition, Geode offers highly optimized shared-nothing disk persistence for in-memory data which does not appear to be available with Ignite.

Apache Cassandra is a highly scalable, distributed key-value store that focuses on eventual consistency. It uses log-structured merge trees to handle writeheavy workloads. The Geode distributed in-memory cluster provides highly performant and advanced capabilities including network partition detection and recovery, version-based state synchronization and conflict resolution, and single IO disk operations.

Apache ActiveMQ and its sub project Apache Apollo offers a powerful message queue framework that is being considered for an open source implementation of Geode's WAN replication capabilities.

Apache Storm is a streaming engine that processes events through a directed graph of computation. It requires Apache Zookeeper for coordination and Apache Kafka to reliably store the streaming data source. Geode provides builtin capabilities for these functions. In addition, Geode offers data locality for related entities and in-process access to reference data typically used during processing.

Apache Kafka offers distributed and durable publish-subscribe messaging. Geode expands beyond publish-subscribe messaging to support data oriented notifications delivered from highly available, fault tolerant event queues that can be easily correlated to related data for further processing. Apache Samza is a distributed processing framework heavily dependent on Apache Kafka and Apache Hadoop YARN for messaging and distributed fault-tolerant processing. Geode is data source agnostic and leverages its own technology and implementation for such use cases.

# Known Risks

Development has been sponsored mostly by a single company (or its predecessors) thus far and coordinated mainly by the core Pivotal GemFire team.

For the project to fully transition to the Apache Way governance model, development must shift towards the meritocracy-centric model of growing a community of contributors balanced with the needs for extreme stability and core implementation coherency.

The tools and development practices in place for the Pivotal GemFire product are compatible with the ASF infrastructure and thus we do not anticipate any on-boarding pains. Migration from the current GitHub repository is also expected to be straightforward.

The project currently includes a modified version of the JGroups software toolkit. JGroups was initially released under the LGPL license. Although we have complied with the terms of the LGPL by making the modified source available, LGPL is classified as an incompatible license by the ASF. The JGroups project has since been re-licenced under ALv2 (see http://www.jgroups.org/license.html) and we plan to engage with the licensor to overcome this license incompatibility. If the license incompatibility cannot be overcome through our discussions with the licensor, we will need to rework this portion of the project based upon a newer version of the JGroups toolkit or via other alternative development efforts.

#### **Orphaned products**

Pivotal is fully committed to Pivotal GemFire and the product will continue to be based on the Geode project. Moreover, Pivotal has a vested interest in making Geode succeed by driving its close integration with sister ASF projects. We expect this to further reduces the risk of orphaning the product.

#### **Inexperience with Open Source**

Pivotal has embraced open source software since its formation by employing contributors/committers and by shepherding open source projects like Cloud Foundry, Spring, RabbitMQ and MADlib. Pivotal also supports other open source projects such as Redis, Chorus, Groovy and Grails. We have experience with the formation of vibrant communities around open technologies with the Cloud Foundry Foundation. Although some of the initial committers have not been developers on an entirely open source, community-driven project, we expect to bring to bear the open development practices that have proven successful on longstanding Pivotal open source projects to the Geode project. Additionally, several ASF veterans agreed to mentor the project and are listed in this proposal. The project will rely on their guidance and collective wisdom to quickly transition the entire team of initial committers towards practicing the Apache Way.

#### **Homogeneous Developers**

While most of the initial committers are employed by Pivotal, we have already seen a healthy level of interest from our existing customers and partners. We intend to convert that interest directly into participation and will be investing in activities to recruit additional committers from other companies.

#### **Reliance on Salaried Developers**

Most of the contributors are paid to work in the Big Data space. While they might wander from their current employers, they are unlikely to venture far from their core expertises and thus will continue to be engaged with the project regardless of their current employers.

#### **Relationships with Other Apache Products**

As mentioned in the Alignment section, Geode may consider various degrees of integration and code exchange with Apache Ignite (incubating), Apache Hadoop, Apache Storm, Apache Spark and Apache Kafka/Samza. Given the success that the Pivotal GemFire product enjoyed as an embedded service, we expect integration points to be inside and outside the project. We look forward to collaborating with these communities as well as other communities under the Apache umbrella.

#### An Excessive Fascination with the Apache Brand

While we intend to leverage the Apache 'branding' when talking to other projects as testament of our project's 'neutrality', we have no plans for making use of Apache brand in press releases nor posting billboards advertising acceptance of Geode into Apache Incubator.

# Documentation

See documentation for the current state of the project documentation available as part of the GitHub repository at https://github.com/project-geode/docs and live at http://geode-docs.cfapps.io/

# **Initial Source**

Pivotal is releasing the source code for Geode under an Evaluation License at https://network.pivotal.io/products/project-geode . We encourage ASF community members interested in this proposal to download the source code, review and try out the software.

# Source and Intellectual Property Submission Plan

As soon as Geode is approved to join Apache Incubator, the source code will be transitioned via the Software Grant Agreement onto ASF infrastructure and in turn made available under the Apache License, version 2.0. We know of no legal encumberments that would inhibit the transfer of source code to the ASF.

# **External Dependencies**

Embedded dependencies (relocated):

- json
- jgroups
- joptsimple

Runtime dependencies:

- antlr
- classmate
- commons-fileupload
- commons-io
- commons-lang
- commons-modeler
   fastutil
- findbugs annotations
- guava
- jackson
- jansi
- javax.activation
- javax.mail-api
- javax.resource-api

- javax.servlet-api
- javax.transaction-api
- ٠ jetty
- ٠ jline
- jna
- json4s log4j
- mx4j
- paranamer
- scala
- slf4j
- snappy-java
- spring
- swagger

Module or optional dependencies:

None

Build only dependencies:

None

Test only dependencies:

- cglib
- hamcrest
- jmock
- junit
- multithreadedtc
- objenesis

Cryptography N/A

# **Required Resources**

#### **Mailing lists**

- private@geode.incubator.apache.org (moderated subscriptions)
- commits@geode.incubator.apache.org
- dev@geode.incubator.apache.org
- issues@geode.incubator.apache.org
- user@geode.incubator.apache.org

#### **Git Repository**

https://git-wip-us.apache.org/repos/asf/incubator-geode.git

#### **Issue Tracking**

JIRA Project Geode (GEODE)

### **Other Resources**

Means of setting up regular builds for Geode on builds.apache.org

# **Initial Committers**

- Amey Barve
- Adib Saikali
- Alan Strait
- · Amogh Shetkar
- Anil GingadeAnilkumar Gingade
- Anthony Baker
- Ashvin Agrawal
- Asif Shahid
- Avinash Dongre
- Barry Oglesby
- Ben Reser
- Bruce Schuchardt
- Bruce Szalwinski
- Catherine Johnson
- Chip Childers

- Christian Tzolov ٠
- Dan Smith
- Darrel Schneider •
- Dave Muirhead
- David Yozie
- Dick Cavender
- Edin Zulich
- · Eric Shu
- Gideon Low
- Greg Chase
- Hemant Bhanawat
- Henry Saputra
- Hitesh Khamesra •
- Jacob Barrett
- Jags Ramnarayan •
- Jan Iversen ٠
- Jason Huynh Jens Deppe ٠
- Jianxia Chen
- •
- John Blum ٠
- Justin Erenkrantz • Ketan Deshpande
- Kirk Lund
- Kishor Bachhav
- Konstantin Boudnik
- Konstantin Ignatyev
- Lise Storc
- Luke Shannon
- Lyndon Adams
- Lynn Gallinat
- Lynn Hughes-Godfrey
- Mark Bretl
- Michael Schubert
- Namrata Thanvi
- Neeraj Kumar
- Nilkanth Patel
- Qihong Chen
- ٠ Rahul Diyewar
- Randy May
- Roman Shaposhnik Severine Tymon •
- •
- Shatarupa Nandi
- Shirish Deshmukh
- •
- Sonal Agarwal Soubhik Chakraborty •
- Sourabh Bansod
- Stephane Maldini Stuart Williams
- •
- Sudhir Menon
- Sunil Jigyasu
- Supriya Pillai
- Suranjan Kumar •
- Suyog Bhokare
- Swapnil Bawaskar •
- Swati Sawant
- Tushar Khairnar
- Udo Kohlmeyer
- Vince Ford
- Vinesh Prasanna Manoharan
- Vivek Bhaskar
- Wes Williams
- William A. Rowe Jr.
- William Markito
- Will Schipp
- Xiaojian Zhou •
- ٠ Yogesh Mahajan

### Affiliations

- WANDisco: Konstantin Boudnik
- Bloomberg LP: Justin Erenkrantz
- Cloud Foundry Foundation: Chip Childers
- NASA JPL: Chris Mattmann
- Unaffiliated: Jan Iversen
- CDK Global: Ben Reser, Konstantin Ignatyev, Bruce Szalwinski ٠
- Pivotal: everyone else on this proposal

# Sponsors

### Champion

Roman Shaposhnik

### **Nominated Mentors**

The initial mentors are listed below:

- Chip Childers Apache Member, Cloud Foundry Foundation
- Justin Erenkrantz Apache Member, Bloomberg LP
  Konstantin Boudnik Apache Member, WANDisco
- Jan Iversen Apache Member, Self employed
- William A. Rowe Jr. Apache Member, Pivotal
- Henry Saputra Apache Member, Pivotal
  Roman Shaposhnik Apache Member, Pivotal
- Chris Mattmann Apache Member, NASA JPL

### **Sponsoring Entity**

We would like to propose Apache incubator to sponsor this project.