

AMELIA

IPsoft's Cognitive Agent

Amelia V3 Platform Design (version 1.4)



This AMELIA® documentation is copyright © 2018 IPsoft Inc and its affiliated companies. All rights reserved.

This document is considered the confidential information of IPsoft and its affiliates. Disclosure to other parties is prohibited unless agreed to in a license or confidentiality agreement

Trademarks, including IPsoft®, AMELIA®, and the IPsoft and AMELIA logos, are the intellectual property of IPsoft Incorporated and its affiliated companies. Any other marks or intellectual property remain the property of their respective licensors or owners.

Table of Contents

1. AMELIA HIGH LEVEL ARCHITECTURE DESIGN	4
1.1 3-NODE CLUSTER ARCHITECTURE	5
1.2 6-NODE CLUSTER ARCHITECTURE	6
1.3 HIGH AVAILABILITY ARCHITECTURE	7
1.4 SINGLE NODE ARCHITECTURE.....	8
2. DEPLOYMENT MODELS	9
2.1 AMELIA HOSTED IN IPSOFT’S CLOUD	9
2.2 CUSTOMER PREMISE / CLOUD.....	10
2.2.1 ONLINE DEPLOYMENT	11
2.2.2 OFFLINE DEPLOYMENT	12
3. HARDWARE/RESOURCE RECOMMENDATIONS	13
3.1 CPU REQUIREMENTS.....	13
3.2 STORAGE REQUIREMENTS.....	14
3.3 DISK IO	14
3.4 CLUSTERED SETUP.....	15
3.5 SINGLE HOST.....	16
3.6 VMWARE OVERSUBSCRIPTION	16
3.7 VMWARE DRS REQUIREMENTS	16
4. COMPONENTS OF AMELIA	17
4.1 OPERATING SYSTEM REQUIREMENTS.....	17
4.2 EXTERNAL LOAD BALANCERS.....	17
4.3 HAPROXY	18
4.4 MYSQL / PERCONA XTRADB CLUSTER (PXC).....	18
4.5 RABBITMQ.....	18
4.6 REDIS SENTINEL	18
4.7 AMELIA DAEMONS	19
4.7.1 AMELIA-ADMIN-WEB	19
4.7.2 USER WEB	19
4.7.3 ENGINE SERVICE.....	19
4.7.4 BATCH SERVICE	19
4.7.5 ESCALATION SERVICE	19
4.7.6 INTEGRATION SERVICE.....	19
4.7.7 DUCKLING SERVICE	20
4.8 TEXT TO SPEECH (TTS).....	20
5. CHAT INTEGRATION.....	21
6. SECURITY	22
6.1 AUTHENTICATION SYSTEMS	22
6.1.1 LOCAL AUTHENTICATION	22
6.1.2 LDAP / ACTIVE DIRECTORY (AD).....	22

6.1.3	DENY ALL	22
6.1.4	SAML 2	22
6.2	ANONYMOUS ACCESS	24
6.3	SSL CERTIFICATES	24
6.4	FIREWALL RULES	24
7.	MONITORING.....	30
8.	BACKUPS.....	32
9.	DISASTER RECOVERY	33
10.	HIGH AVAILABILITY CONFIGURATIONS.....	34
10.1	FRONT END SERVICES	35
10.2	ADMIN POD	35
10.3	RABBITMQ HOSTS.....	35
10.4	CONVERSATION POD.....	36
10.5	GENERAL NOTES ABOUT HARDWARE/RESOURCE RECOMMENDATIONS.....	36
10.5.1	CPU.....	36
10.5.2	DISK.....	36
10.5.3	VMWARE.....	36
10.5.4	OS.....	37
10.6	SCALING AMELIA COMPONENTS	37
10.6.1	DATABASES	37
10.6.2	SCALING ADB AND KDB.....	37
10.6.3	SCALING CDB	37
10.6.4	ALTERNATE CDB DESIGN (NOT RECOMMENDED)	37
10.6.5	RABBIT MQ.....	38
10.6.6	INTEGRATION SERVICES	38

Figures

Figure 1.	Service Architecture Diagram.....	4
Figure 2.	Initial 3-Node Cluster Environment Diagram	5
Figure 3.	6-Node Cluster Environment Diagram	6
Figure 4.	Amelia Architecture Scaled for 1200 Concurrent Conversations	7
Figure 5.	1-Node Cluster Environment Diagram	8
Figure 6.	Amelia Hosted in IPsoft's Cloud with IPsec VPN Tunnels	9
Figure 7.	Amelia Hosted in IPsoft's Cloud	10
Figure 8.	Amelia Hosted in Customer's Cloud	11
Figure 9.	Chat Integration	21
Figure 10.	Amelia and Service Provider Initiated Authentication.....	23
Figure 11.	Amelia and Identity Provider Initiated Authentication.....	23
Figure 12.	Example of OS Monitoring Checks Performed	30

Figure 13. Amelia Deployment Scaled for 1200 Concurrent Conversations..... 34

Tables

Table 1. Preferred Processors	14
Table 2. 3-Node Clustered Setup Environment Resource Requirements	15
Table 3. 6-Node Clustered Production Environment Resource Requirements	15
Table 4. Conversation POD Setup Environment Resource Requirements	16
Table 5. Single Host Environment Resource Requirements.....	16
Table 6. SSL Certificate Requirements	24
Table 7. Firewall Rules.....	24

Document History

Author	Version	Date	Comments	Final Approval?
Randy Schneiderman	0.1	11/15/2017	Initial Document	Yes
Randy Schneiderman	1.0	11/30/2017	Updated for V3	Yes
Randy Schneiderman	1.1	1/18/2018	Added NetData and DR sections	Yes
Randy Schneiderman	1.2	1/29/2018	Revised Sizing Requirements	Yes
Randy Scheiderman	1.3	6/21/2018	Add Duckling Service and POD details, updated Figure 1	Yes
IPsoft Research & Development	1.4	12/10/2019	Added high availability section	Yes

1. Amelia High Level Architecture Design

Amelia is the artificial intelligence platform that can understand, learn and interact as a human would to solve problems. Amelia reads natural language, understands context, applies logic, infers implications, learns through experience and even senses emotions. The diagrams below illustrate Amelia's architecture.

Amelia V3 is designed with multiple shared services, including but not limiting to Administrative Services, Conversation Engine Pods, Integration Services, and a Database Shard Architecture. The concept for this design is to separate very large JVMs, databases, daemons into smaller, faster, more easily managed fragments. The below diagrams show the service architecture and data flow overviews of Amelia V3:

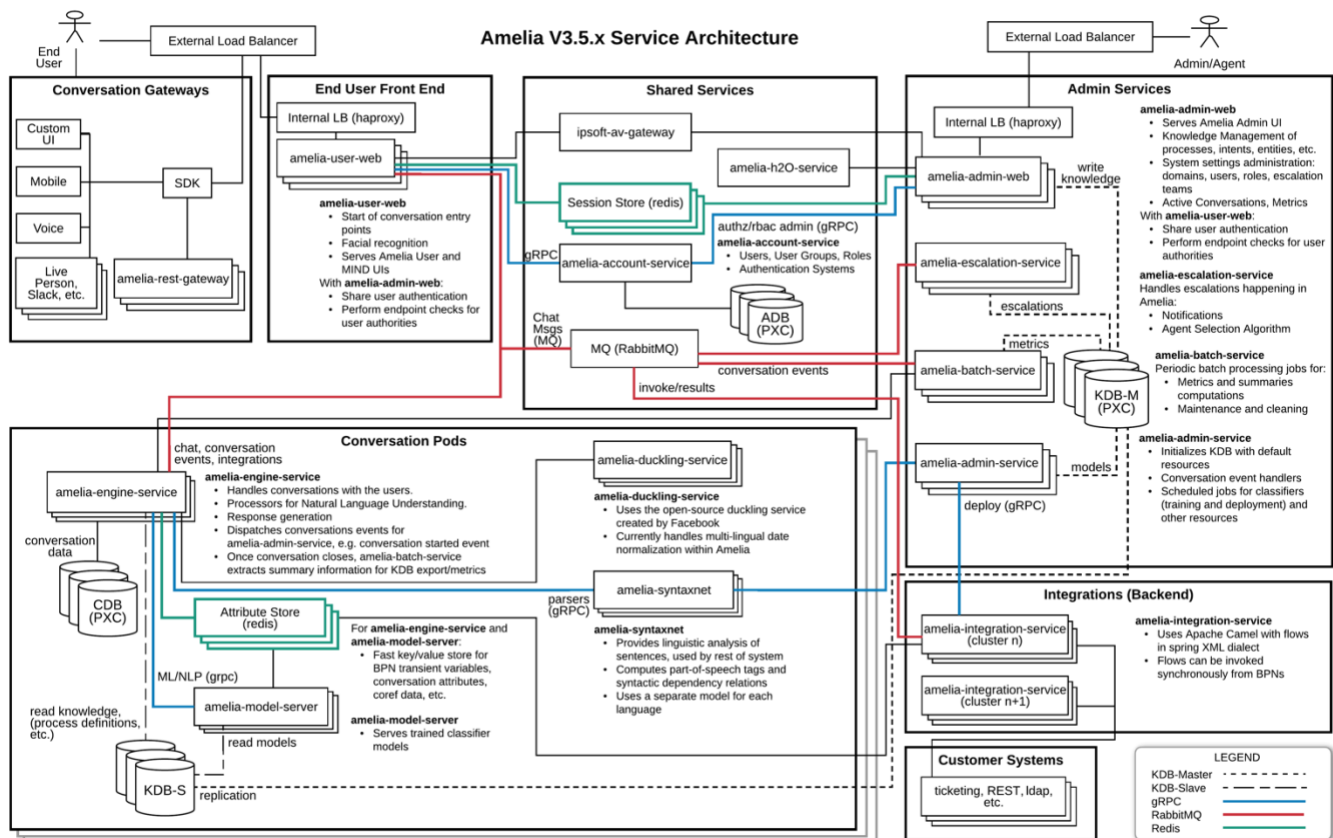


Figure 1. Service Architecture Diagram

The below sections will describe the infrastructure requirements and necessary components of Amelia V3. Some hosts will have multiple middleware components to support Amelia and her functions. Amelia supports clustering both for scaling and high availability. It is recommended that all production environments be clustered with three Application servers and three Database servers. When configured with clustering, there are no single points of failure within Amelia; any failure of a single component should at most result in a few

seconds of intermittent faults. Amelia V3 is only supported on LANs and where network split-brain between nodes is very unlikely. In the event of a network split-brain event, manual intervention may be required and data in-flight may be lost.

1.1 3-NODE CLUSTER ARCHITECTURE

IPsoft's standard deployment model, a best practice, consists of three identical servers (see Figure 1).. When clustered, there are no single points of failure within Amelia and a failure of any single component should at most incur a few seconds of intermittent errors.

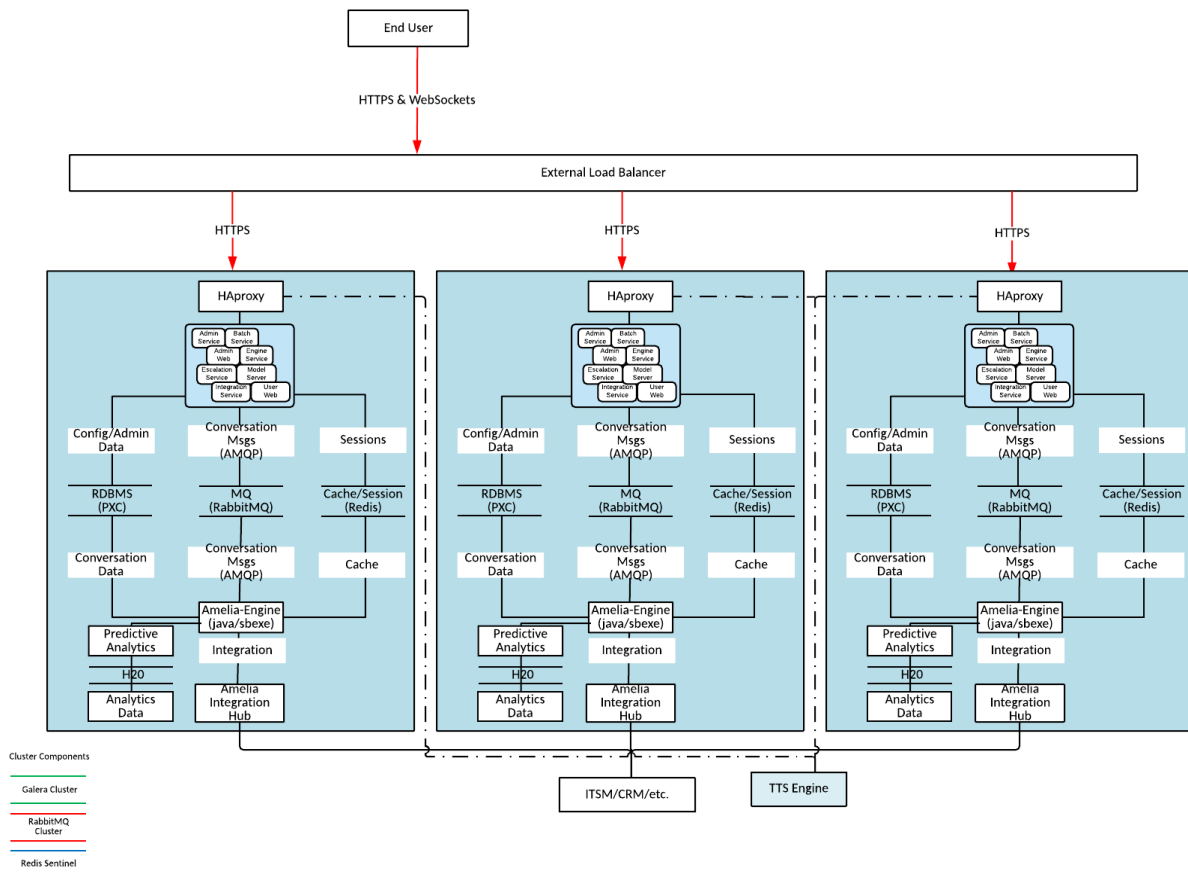


Figure 2. Initial 3-Node Cluster Environment Diagram

NOTE

Connection options — for example, to IPcenter, LDAP/SSO/SAML2, or other technology — should be discussed with IPsoft in the planning process.

1.2 6-NODE CLUSTER ARCHITECTURE

This configuration splits the application and database services into two network tiers. As Amelia V3 uses a sharding/multiple shared services architecture, the various application and database services can be further split off into additional servers to provide separation of service requirements and for large volume use cases.

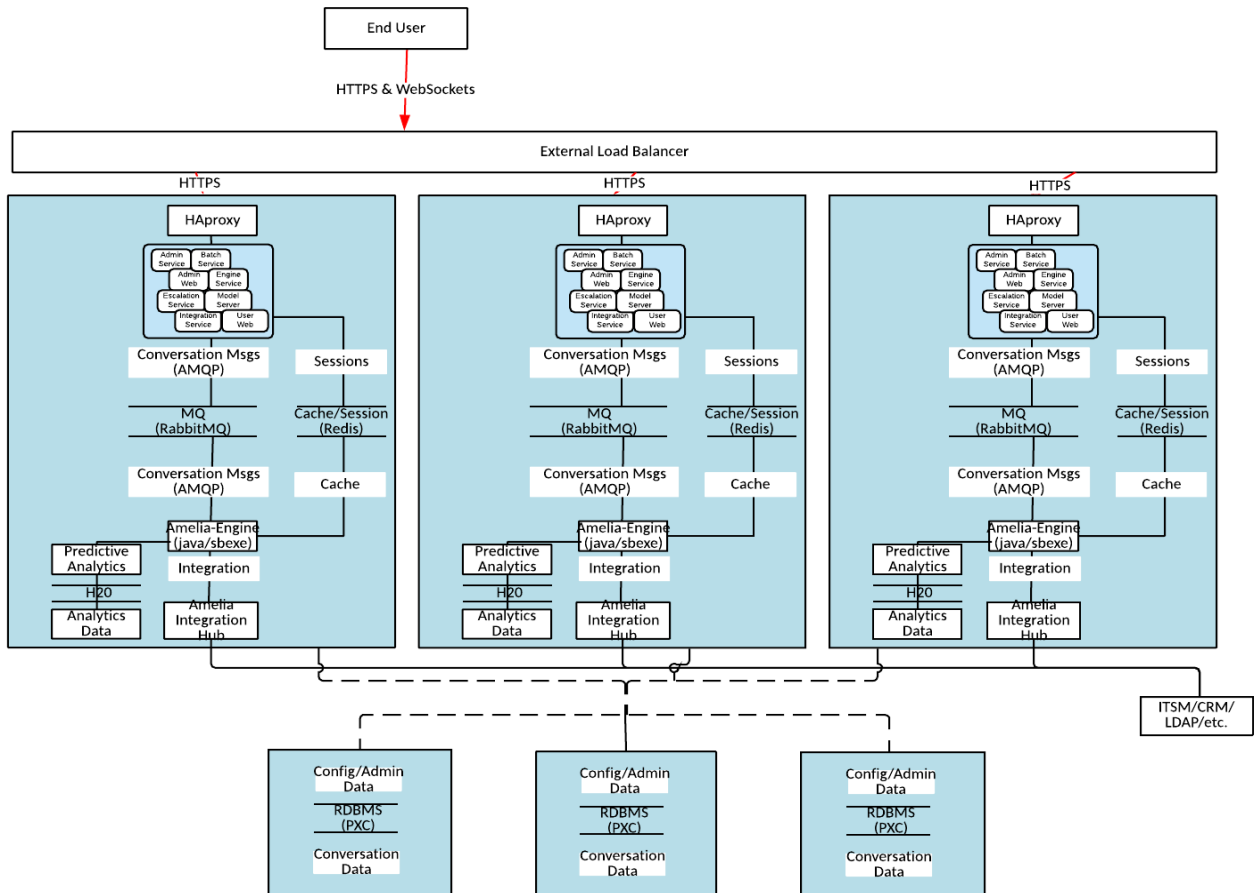


Figure 3. 6-Node Cluster Environment Diagram

NOTE

Connection options — for example, to IPcenter, LDAP/SSO/SAML2, or other technology — should be discussed with IPsoft in the planning process.

1.3 HIGH AVAILABILITY ARCHITECTURE

Refer to the High Availability Configuration section below for more details about this configuration.

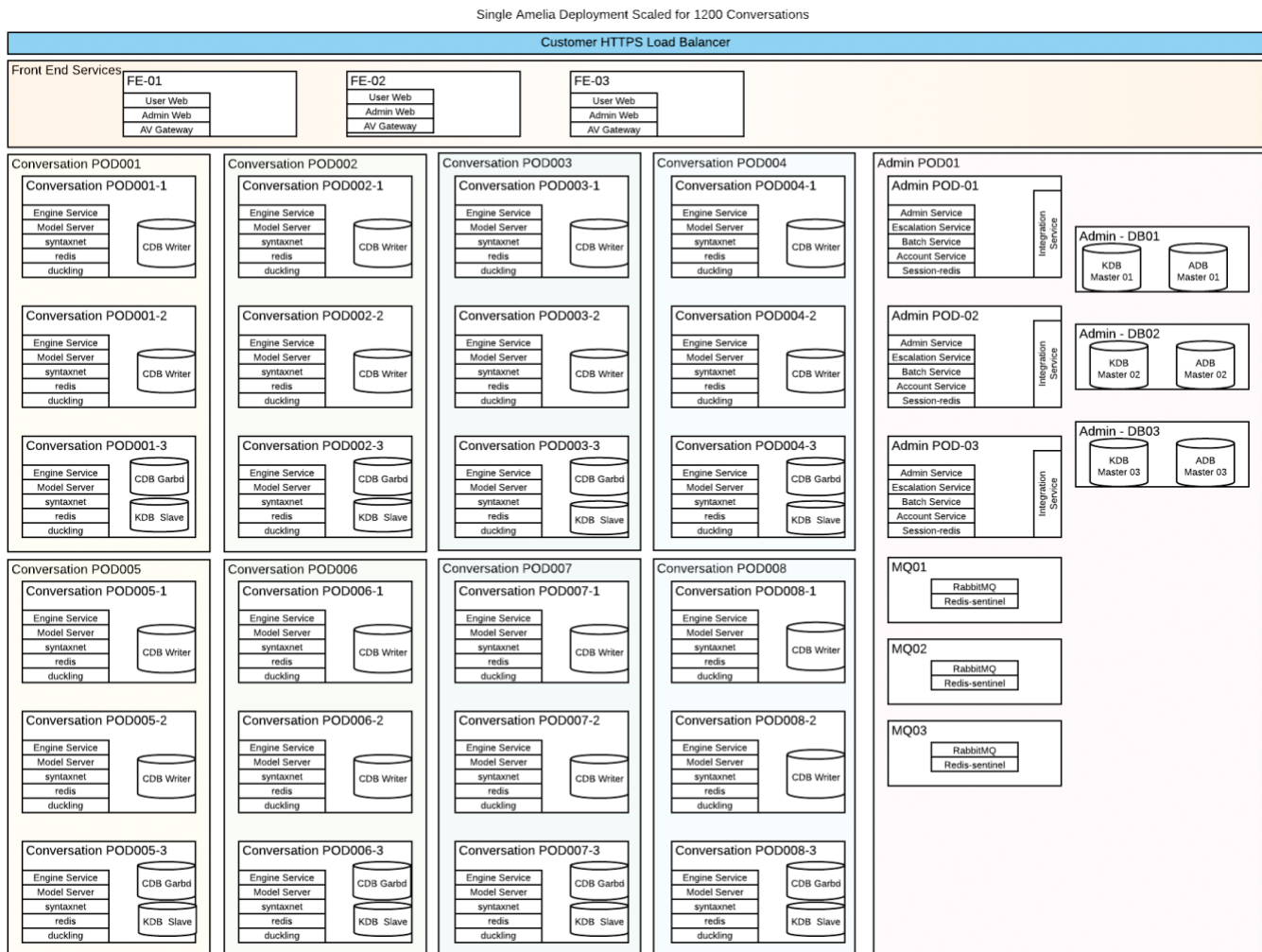


Figure 4. Amelia Architecture Scaled for 1200 Concurrent Conversations

NOTE

Connection options — for example, to IPcenter, LDAP/SSO/SAML2, or other technology — should be discussed with IPSOFT in the planning process.

1.4 SINGLE NODE ARCHITECTURE

Amelia also supports single host architectures. Single-host configurations are meant for POC/Dev environments.

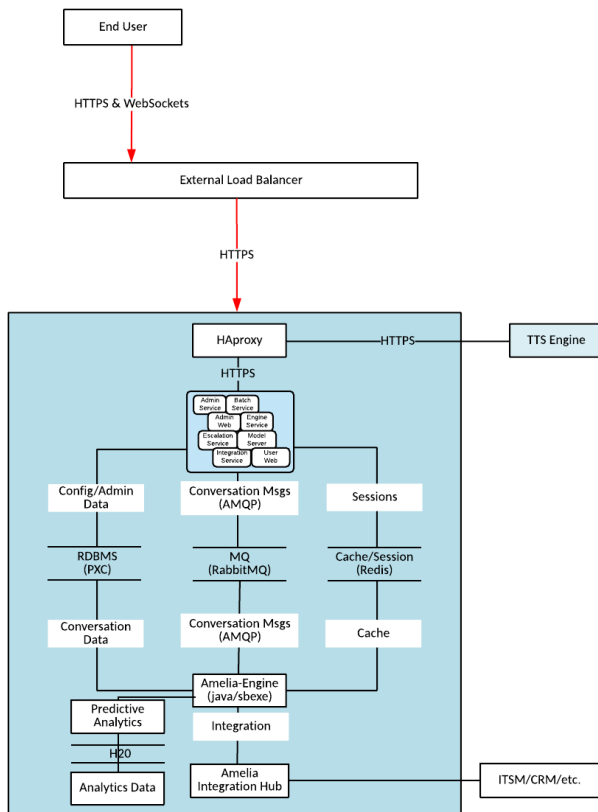


Figure 5. 1-Node Cluster Environment Diagram

NOTE

Connection options — for example, to IPcenter, LDAP/SSO/SAML2, or other technology — should be discussed with IPsoft in the planning process.

2. Deployment Models

Amelia can be delivered in several deployment models, depending on customer/partner requirements. On-Premise deployments will may require encrypted network connectivity between IPsoft and the clients/partners environments for installation of Amelia V3. Interconnectivity can be accomplished with dedicated communication circuits and/or site to site Virtual Private Network (“VPN”) tunnels across the public Internet.

IPsoft has developed a self-contained tool called Amelia Deployment Center (ADC) for deployment and configuration of Amelia V3 for remote servers to install Amelia on RHEL 7 family servers, which is used to automate and organize system configuration tasks.

2.1 AMELIA HOSTED IN IPSOFT’S CLOUD

Amelia can be hosted at IPsoft’s datacenters worldwide in IPsoft’s New York Metro and Amsterdam datacenters. In this model, IPsoft assumes all responsibility for deployment and scaling of Amelia for given clients and partners. This deployment will utilize IPsoft’s current hardware comprised of Dell servers, Compellent storage, VMware, and Cisco/Arista networking.

For security purposes, backend technologies and any integrations with Amelia are interconnected utilizing a secure delivery network, typically an IPsec VPN and/or MPLS.

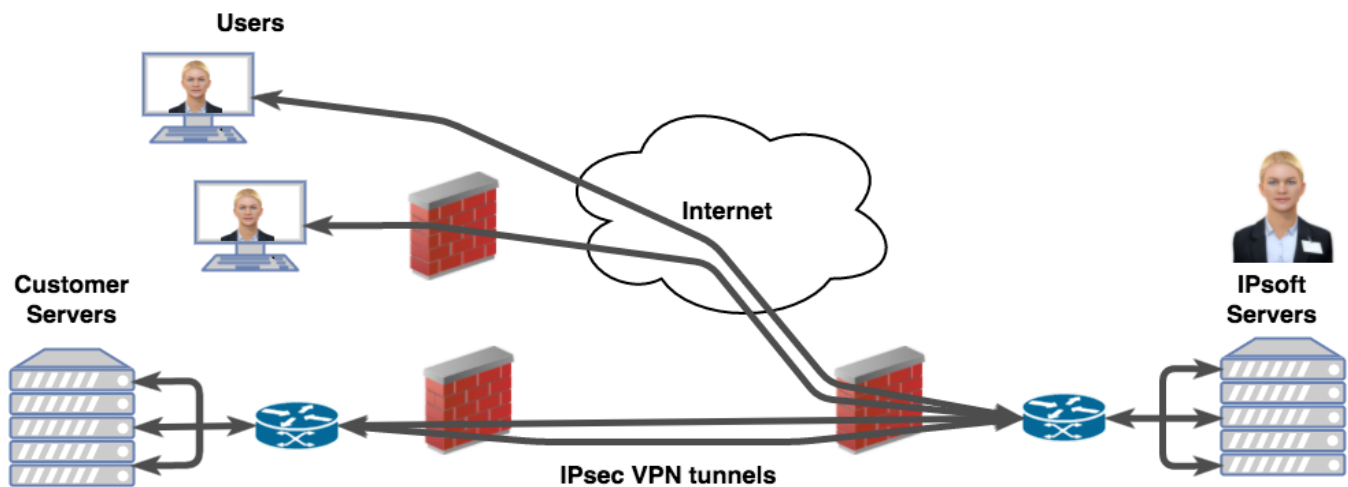


Figure 6. Amelia Hosted in IPsoft’s Cloud with IPsec VPN Tunnels

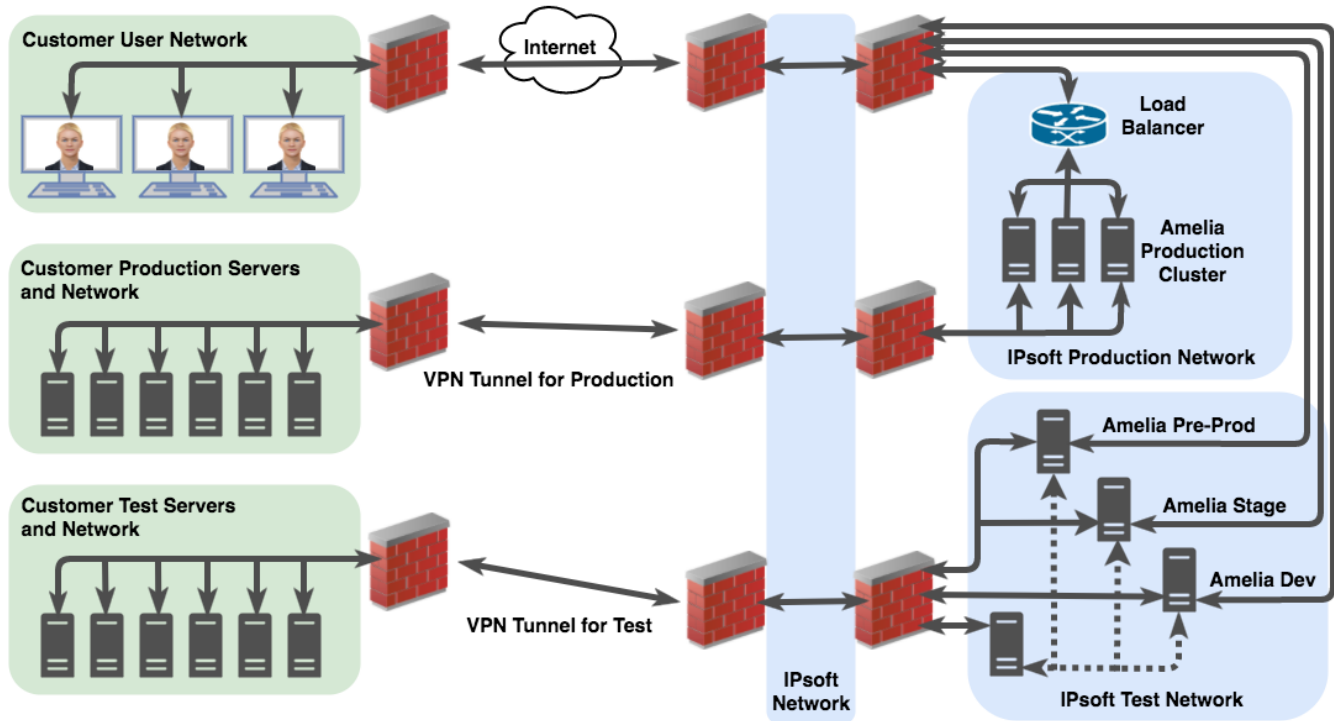


Figure 7. Amelia Hosted in IPsoft's Cloud

Figure 6 shows a conceptual view of how clients and servers are communicating with Amelia in IPsoft-hosted environment.

2.2 CUSTOMER PREMISE / CLOUD

Amelia can also be deployed within a client's/partner's facilities, without IPsoft managing the hardware and network infrastructure. This type of deployment requires the involvement of Client/Partner Architects and IPsoft's Service Design resources to develop a joint architecture for connectivity and sizing.

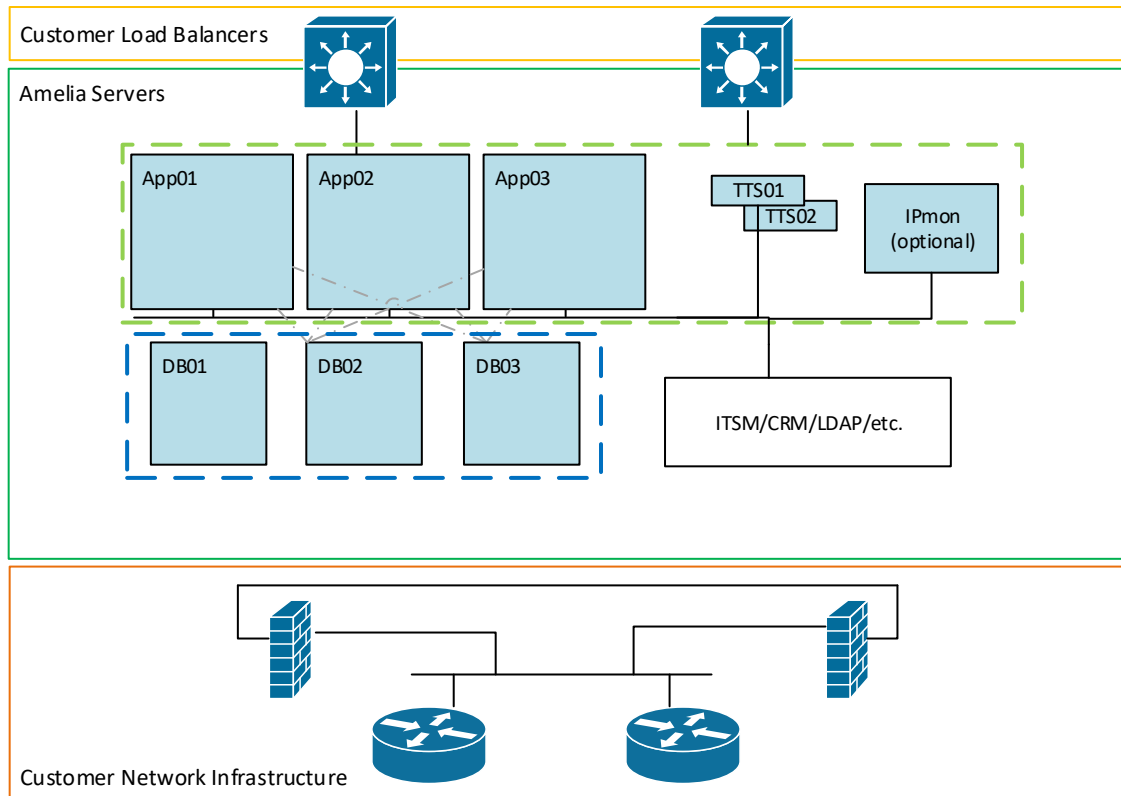


Figure 8. Amelia Hosted in Customer's Cloud

NOTE

Connection options — for example, to IPcenter, LDAP LDAP/SSO/SAML2, or other technology — should be discussed with IPsoft in the planning process.

2.2.1 Online Deployment

With remote connectivity enabled, IPsoft will use ADC to deploy Amelia to automate the installation of Amelia and her dependencies. Please refer to IPsoft's Amelia Deployment Center (ADC) Guide for in-depth information.

The deployment and release of Amelia to remote server(s) is managed by IPsoft's implementation of ADC. Using job templates and ADC mechanism, IPsoft can show progress and monitor the status of Amelia being deployed.

Deployments in the public cloud (such as AWS/Azure) are considered on-premise deployments, as the roles and responsibility of the infrastructure is owned by the client/partner. IPsoft has done extensive testing within AWS for all support Operating Systems as described below in Section 3.1

At a minimum, each server should be deployed using the r4.xlarge sizing; this does not include a backup partition if that is required. Client is responsible for any Elastic Load Balancing configuration and any OS/network security aspects.

As of this writing, Amelia is not supported using containers, for example, Docker.

2.2.2 Offline Deployment

At times, remote connectivity will not be possible between IPsoft and client/partner for deployment of Amelia on remote server(s). IPsoft refers the lack of a persistent connection as an Offline Instances. IPsoft will continue to use ADC for deployments and the initialization of all playbooks will be maintained on the server(s) locally within the client's/partner's network.

For offline deployments, IPsoft can only provide specific SLAs based on available connectivity and access.

If considering an Offline deployment of Amelia, topics of discussion should include but are not limited to:

- Backups – How are they performed? Schedule?
- Monitoring – OS, Middleware, Amelia, Web Checks
- Administration – Engineer access to Operating System
- Upgrades – OS, Amelia
- Support – Break/Fix, Root level access
- Deployment Timeline– Procurement of servers/load balancers and/or access limitations

For offline deployments with a pre-existing IPcenter instance, contact IPsoft to discuss options.

3. Hardware/Resource Recommendations

For deployments of Amelia on premise /cloud, there are several recommendations and requirements, depending on the type(s) of use cases and volume for each use case. Amelia can be scaled up/out by increasing the number of CPUs, RAM, and disk capacity as well the number of Conversation PODs . If Amelia is configured as a single node, It cannot be scaled-out with additional nodes.

Amelia deployments are supported in both physical and virtual environments. For virtual deployments, it is recommended to enable Memory and Disk reservations for best performance; especially in highly utilized shared infrastructure.

3.1 CPU REQUIREMENTS

Amelia does not perform well in virtual infrastructures with Memory ballooning and/or CPU scheduling issues. Virtual machines (VMs) depend on available host resources (CPU, Memory), and the guest operating system consumes those resources. A problem with resource availability or scheduling inside or outside the virtual machine may cause it to become unresponsive.

Reviewing CPU performance metrics can be used to determine whether a guest operating system is actually running, whether the virtual machine monitor (VMM) is running, or whether there is scheduling contention. The metrics is also leveraging insight into the responsiveness of a virtual machine or its Guest OS:

- Run - Amount of time the virtual machine is consuming CPU resources.
- Wait - Amount of time the virtual machine is waiting for a VMkernel resource.
- Ready - Amount of time the virtual machine was ready to run, waiting in a queue to be scheduled.
- Co-Stop - Amount of time a SMP virtual machine was ready to run, but incurred delay due to co-vCPU scheduling contention.

It is recommended Memory/CPU Hotplug is enabled to increase CPUs/RAM at any time to scale quickly. Be mindful of the hardware configuration and NUMA settings, it is optimal for Amelia to access CPU/Memory resources on the same NUMA node.

Amelia should be utilizing at a minimum dual Intel® Xeon® Processor E5-2687W v3 (25M Cache, 3.10 GHz).

These are preferred processors, from newest to oldest.

Table 1. Preferred Processors

Product Line/Processor Name	Cores	Threads	Base Frequency	Max Turbo Frequency	Cache
Intel 2nd Gen Scalable Processors					
Intel® Xeon® Gold 6246 Processor	12	24	3.30 GHz	4.20 GHz	24.75 MB
Intel® Xeon® Gold 6242 Processor	16	32	2.80 GHz	3.90 GHz	22 MB
Intel® Xeon® Platinum 8268 Processor	24	48	2.90 GHz	3.90 GHz	35.75 MB
Intel Scalable Processors					
Intel® Xeon® Gold 6136 Processor	12	24	3.00 GHz	3.70 GHz	24.75 MB L3
Intel® Xeon® Gold 6142 Processor	16	32	2.60 GHz	3.70 GHz	22 MB L3
Intel® Xeon® Platinum 8158 Processor	12	24	3.00 GHz	3.70 GHz	24.75 MB L3
Intel E5 V4 Processors					
Intel® Xeon® Processor E5-2687W v4	12	24	3.00 GHz	3.50 GHz	30 MB SmartCache

3.2 STORAGE REQUIREMENTS

Amelia's response time (latency) is important for the user's experience, handling of high concurrent connections, and scalability. Amelia requires a low latency and high throughput storage tier. IPsoft highly recommends using Solid State Drives (SSDs) for primary storage for all Amelia Servers. SAS/SATA Storage can be used for backups if desired.

Customers and Partners can leverage existing storage platforms that are installed with SSDs. Both All-Flash-Arrays and hybrid (SSDs and Spinning Disks) are viable solutions, taking into account the initial write lands on SSD storage tier.

Amelia's databases and PODs require at minimum 10,000 IOPS, 15,000 to 20,000 for larger deployments.

3.3 DISK IO

Each VM must be able to sustain 1000 Read and 1000 Writes IOPS at <10ms latency.

3.4 CLUSTERED SETUP

For Production/DR deployments, it is highly recommended to deploy a clustered setup for high availability and scaling for performance, Amelia can be deployed on physical or virtual servers. As of this writing, the following are guidelines on the infrastructure requirements for a clustered Production/DR deployment. These requirements are for each server.

Table 2. 3-Node Clustered Setup Environment Resource Requirements

Tier	Specifications per Host			Number of Hosts
	CPU*	RAM (GB)**	Disk Capacity (GB)***	
Amelia Node	16	80	512 (not including OS)	3
TOTALS	48	240	1.6TB	3

*Based on underlying server's NUMA settings, multiple sockets may be suggested

**Additional RAM/Disk Capacity for language pack support & gateways (1GB RAM / 10GB disk capacity for each)

***/apps is the required mount point; LVM and/or XFS file system is highly recommended. Slower disks can be utilized if slow performance is acceptable. Please see the below notes regarding proper disk capacity sizing.

NOTE

The above table are requirements for each node of the cluster.

Table 3. 6-Node Clustered Production Environment Resource Requirements

Tier	Specifications per Host			Number of Hosts
	CPU*	RAM (GB)**	Disk Capacity (GB)***	
Apps	12	64	300	3
Database	8	24	1024 (not including OS)	3
TOTALS	60	264	3.88TB	6

*Based on underlying server's NUMA settings, multiple sockets may be suggested

**Additional RAM/Disk Capacity for language pack support & gateways (1GB RAM / 10GB disk capacity for each)

***/apps is the required mount point; LVM and/or XFS file system is highly recommended. Slower disks can be utilized if slow performance is acceptable. Please see the below notes regarding proper disk capacity sizing.

NOTE

Clustering is only supported on LANs and where network split-brain between nodes is very unlikely. In the event of a network split-brain event, manual intervention may be required and data in-flight may be lost.

Databases sizes for both 3-node and 6-node clusters depends on the following:

- Data Retention Requirements/Compliance
- Local Database Backup(s)
- Use Cases
- Number of Conversation Pods

Conversation Pods are used for additional capacity for higher concurrent conversations. Conversation Pods are deployed in batches of 3 VMs to handle an additional 150 concurrent conversations. There currently is no limit

on the number of Pods that can be deployed with Amelia. For deployments with high peaks, it is recommended to configure Amelia with the necessary sizing for those peaks and scale down if desired.

Table 4. Conversation POD Setup Environment Resource Requirements

Tier	Specifications per Host			Number of Hosts
	CPU*	RAM (GB)**	Disk Capacity (GB)***	
POD Node	8	32	100 (not including OS)	3
TOTALS	24	96	300 GB	3

3.5 SINGLE HOST

For deployments of Proof-of-Concepts, Development, User Acceptance Testing (UAT), and some limited Production environments, Amelia can be deployed on a single physical or virtual server. As of this writing, the following are guidelines on the infrastructure requirements for a single host deployment

Table 5. Single Host Environment Resource Requirements

Resource	Quantity	Notes
CPU	16	Based on underlying server's NUMA settings, multiple sockets may be suggested
RAM**	80 GB	
Disk Capacity	300 GB	/apps is the required mount; LVM and/or XFS file system is highly recommended. Slower disks can be utilized if slow performance is acceptable.

**Additional RAM/Disk Capacity for language packs & gateways (1GB RAM/10GB Disk Capacity each pack)

3.6 VMWARE OVERSUBSCRIPTION

Do not oversubscribe CPUs and memory. One vCPU equals one Hyperthread.

3.7 VMWARE DRS REQUIREMENTS

VMware DRS (Distributed Resource Scheduler) is a utility that balances computing workloads with available resources in a virtualized environment. With VMware DRS, users can define rules for the allocation of physical resources and placement among virtual machines. The utility can be configured for manual or automatic control. Resource pools can be easily added, removed or reorganized. If desired, resource pools can be isolated between different business units. If the workload on one or more virtual machines drastically changes, VMware DRS redistributes the virtual machines among the physical servers.

4. Components of Amelia

An Amelia instance includes database, message transport, load balancers, and other components, as well as escalation and other processes and systems to set up.

4.1 OPERATING SYSTEM REQUIREMENTS

Amelia is a Linux based Application and can be deployed onto a RHEL 7-family OS. These are:

- Scientific Linux 7.x (SL)
 - Used for IPsoft Cloud and On-Premise Online Deployments
 - IPsoft is responsible for Amelia code and OS patches
 - IPsoft's images follows the Center for Internet Security® (<http://www.cisecurity.org/>) benchmarks, referred to as CIS.
- CentOS 7.x
 - Licensing is provided via GPL
 - Yum repository for software dependencies is provided by client or via Internet
 - IPsoft's images follows the Center for Internet Security® (<http://www.cisecurity.org/>) benchmarks, referred to as CIS.
- RedHat Enterprise Linux 7.x (RHEL)
 - Licensing is provided by client at time of install
 - Yum repository for software dependencies is provided by client at time of install
 - IPsoft can provide RHEL CIS image if desired
- Oracle Linux 7.x (RHEL)
 - Licensing is provided via GNU General Public License (GPLv2). Support contracts are available from Oracle.
 - Yum repository for software dependencies is provided by client or via Internet
 - IPsoft does not provide a Oracle Linux image

IPsoft can provide an Open Virtualization Appliance (OVA) with the requisite OS, software dependencies, and Amelia herself. Leveraging IPsoft's OVA provides a quick approach to standing up an Amelia instance. On request, IPsoft also can deploy the software dependencies and Amelia on a Client's OS build. Client is responsible for any licensing and yum repository access.

4.2 EXTERNAL LOAD BALANCERS

Amelia clusters (can also be configured for single node deployments) leverage external load balancers to handle layer 4 (transport layer) traffic. Load balancing this way will forward user/API traffic based on host and port availability. Health Checks determines if a backend server is available to process requests. The default health check is to establish a TCP connection to the server on the configured hostname/IP and port.

IPsoft recommends using the “least connections” algorithm because of the potential for longer sessions. Load balancers will forward https connections to the Amelia-Web services; any SSL certificates will be provided to IPsoft to decrypt the SSL traffic.

4.3 HAPROXY

HAProxy is normally used to handle external load balancing requests for web/application loads. IPsoft uses HAProxy for service checks and load balancing internal Amelia services and dependencies. This allows for better utilization of all servers and ensures availability. IPsoft configures HAProxy by binding the loopback IP address (127.0.0.x) as the frontend listening IP/port, forwards traffic via the “Round Robin” algorithm to the other servers, with the necessary health check.

4.4 MYSQL / PERCONA XTRADB CLUSTER (PXC)

The Knowledge Database (KDB) to store configuration parameters such as domain, authentication, FAQ, Grammar, transactions from SLUs/BPNs, classifiers. In addition, there is at least one slave of this database which is used at conversation time to query this information. Writes to this database are only done through the admin daemons.

The CDB (Conversation Database) is used for storing per-conversation data. There may be multiple instances of the CDB database attached to multiple engines. To support additional conversation throughput, additional instances can be added along with the associated engines. These databases are not accessible from the admin daemons.

For more information regarding PXC, here is a web link to Percona: <https://www.percona.com/software/mysql-database/percona-xtradb-cluster>

4.5 RABBITMQ

RabbitMQ is used for various messaging tasks and share data between the Amelia Engines and Web services as a three-node active/active/active cluster. It uses both Stomp and AMQP messaging protocols which are exposed on ports 13351 through 13354. Stomp and AMQP frontends are configured in HAProxy to distribute connections in RabbitMQ. Upon failure of a single node, the engine and web will reconnect and begin consuming and sending messages from one of the remaining nodes.

For more information regarding Rabbit MQ Clustering, here is a web link to RabbitMQ’s documentation: <http://www.rabbitmq.com/clustering.html>

4.6 REDIS SENTINEL

Redis is a Key-Value In-Memory data structure store used for caching and configured with a single master and multiple slaves. Monitoring and automatic master promotion is handled by Redis Sentinel. Amelia connects to a Sentinel front-end in their local HAProxy to retrieve the current master Redis server information and then connect directly to that instance. Should the master Redis instance fail, a new master is elected and failover occurs automatically.

For more information regarding Redis Sentinel, here is a web link to Redis' documentation:

<http://redis.io/topics/sentinel>

4.7 AMELIA DAEMONS

4.7.1 Amelia-Admin-Web

Amelia-Web has a Spring Security layer to authenticate and authorize external connections to the system. Contains the REST APIs used by the UI and has read/write access to the KDB master database; no access to CDB.

4.7.2 User Web

The end user and agent conversation interface. Has read-only to a KDB slave and read/write to a CDB.

4.7.3 Engine Service

Amelia Engine Packs (AEP) are bundled units of language, process, and server-side integrations to achieve conceptual objectives.

4.7.4 Batch Service

The Batch Service provides metric calculation and other batch jobs. Has read/write to the KDB master and no access to CDB.

4.7.5 Escalation Service

Amelia can escalate during a conversation in the following manners:

- Misunderstanding: Core dialog manager or a subsystem is unable to handle the user's response.
- Explicit Escalation: Can occur only from BPN, through an "escalate" task.
- Warm Handover: A special explicit escalation, where a reason is specified.

The Escalation Service has read/write to the KDB master and no access to CDB.

4.7.6 Integration Service

The Integration Service is a process run separately from Amelia, potentially on another host or hosts, to allow Amelia to interface with external systems. Integration Flows are Apache Camel contexts created in Amelia V3 admin tools and deployed to these remote processes over GRPC. The Integration Service unpacks the bundle and deploys it in its own Spring application context separate from that of the parent context and of any other flows running on the Integration Service. Integration Service relies on Spring Integration to handle Direct RPC-Style calls over RabbitMQ, and then internally hands off to Apache Camel to execute the given request.

Integration with Amelia is complex and can vary depending on the usage. Please reach out to your IPsoft's Cognitive Lead for best practices involving integration with specific technologies.

4.7.7 Duckling Service

Based on Facebook's open source Duckling project, this service is implemented it within Amelia to handle multi-lingual date normalization. Duckling service uses probabilistic context free grammar based on rules consisting of patterns (regular expressions for character level and predicates for concept level matching) and productions/derivations. The modules that parse temporal expressions in English, Spanish, French, Italian and Chinese.

4.8 TEXT TO SPEECH (TTS)

TTS may be used for speech synthesis in Amelia. Multiple TTS voice engines may be installed to facilitate synthesis of audio to satisfy specified language requirements. The TTS server receives encoding requests and also serves the front-end content.

The TTS runs on a Windows 2012 R2 server, ideally with a load balanced deployment of two or more Windows Servers in Active/Passive mode behind a provided Virtual IP (VIP) address. IPsoft recommends deploying a TTS server with a minimum of 8CPUs/16GB RAM/100GB disk capacity. Please inquire with RND on supported languages and licensing information.

TTS is not provided as part of Amelia Deployments; this is a separate, licensed installation.

5. Chat Integration

Amelia can be integrated with external chat systems like LivePerson/LiveEngage, Skype for Business, Facebook Messenger. Please reach out to your IPsoft's Cognitive Lead for best practices involving integration with specific technologies.

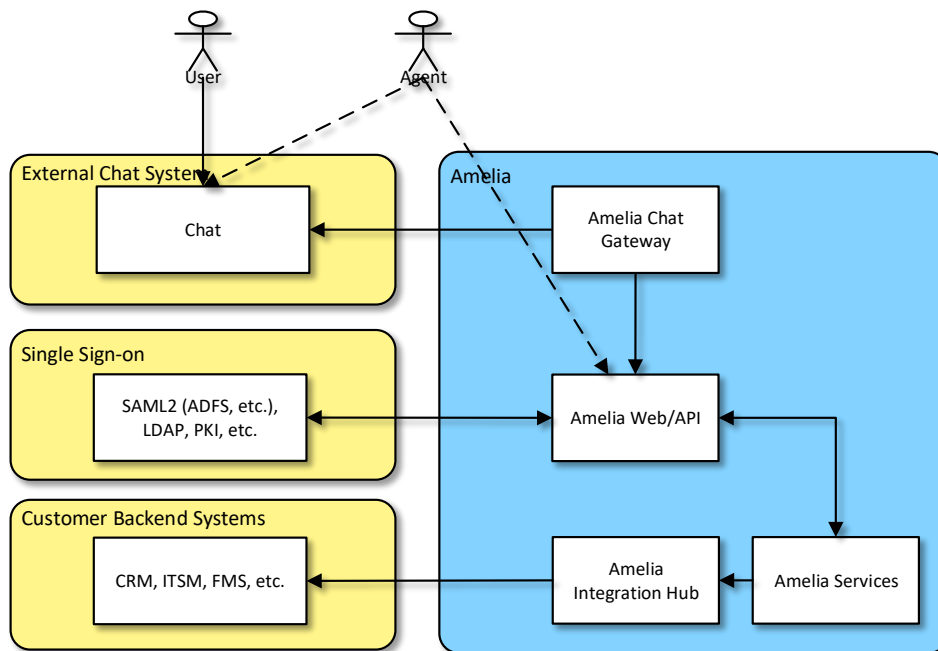


Figure 9. Chat Integration

6. Security

6.1 AUTHENTICATION SYSTEMS

Authentication systems in Amelia define where the user credentials are stored, as well as how to verify credentials. All authentication systems have the following configurable properties. Specific, or custom authentication systems, may require additional configuration or custom development.

6.1.1 Local Authentication

Passwords are stored hashed with Bcrypt in the Amelia database. Local Authentication is the default setting “out-of-the-box”.

6.1.2 LDAP / Active Directory (AD)

For deployments leveraging LDAP, Amelia does not store user passwords (or hashes), but instead delegates password verification to the configured LDAP server or Windows Domain Controllers. The following additional configuration parameters are required to configure LDAP/AD properly.

6.1.3 Deny All

The Deny All authentication system is used for required accounts that should not allow interactive login. Any attempt to authenticate to this authentication system will fail and no passwords or hashes are stored. This authentication system is used for the Amelia user by default.

6.1.4 SAML 2

Amelia supports both Service Provider (SP) initiated and Identity Provider (IDP) initiated SAML2 authentication. As with all authentication methods supported by Amelia, SAML2 is integrated within the core Amelia authentication and authorization framework. In the SAML2 case, support is provided by Spring Security SAML. SAML authentication is only supported for web clients (not mobile).

- Service Provider (SP): An application providing a service to an end-user. In this case, Amelia.
- Identity Provider (IDP): An application/service that manages user identities and provides authentication capabilities. For example, Microsoft Active Directory Federation Services (ADFS) or Ping Federate.

SERVICE PROVIDER INITIATED AUTHENTICATION

In SP initiated authentication, an end user first requests a resource within Amelia. If the resource requires authentication and the user is not yet authenticated with Amelia, the user will be sent to the IDP for authentication. Upon successful authentication, the IDP will send the user back to Amelia with a cryptographically secure assertion of their identity.

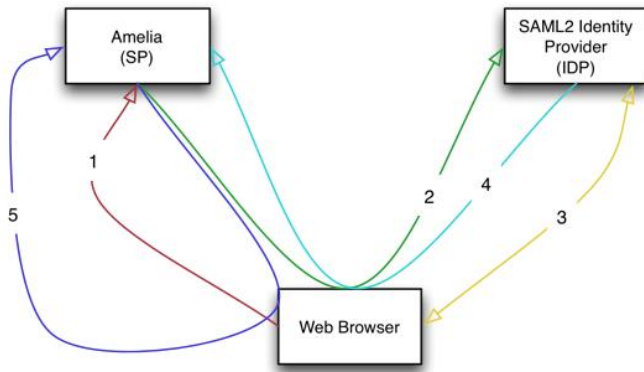


Figure 10. Amelia and Service Provider Initiated Authentication

1. An unauthenticated user attempts to access a protected resource in Amelia.
2. Amelia generates an authentication request (signed AuthnRequest) and redirects the user's browser to the IDP with this request.
3. The IDP determines the user's identity. This could involve one-time tokens, username and password, or other authentication methods. The exact mechanism is not important to the integration.
4. Once the user is authenticated, the IDP generates a signed AuthnResponse carrying the user's identity, email, and other profile information as required. The IDP then redirects the user's browser back to Amelia with this response.
5. Amelia verifies the signed response, and if valid, auto-creates (if required/configured) the user, and logs them into Amelia. Amelia then redirects the user's browser to the original protected resource they requested in step 1.

All communication takes place over TLS and is between the user's browser and the SP, and the user's browser and the IDP. No direct communication is done between the SP and IDP other than initial out of band sharing of metadata at configuration time (see below).

IDENTITY PROVIDER INITIATED AUTHENTICATION

Amelia supports IDP initiated authentication, however, SP initiated authentication should be preferred. In the IDP initiated scenario users must first authenticate to the IDP before accessing Amelia.

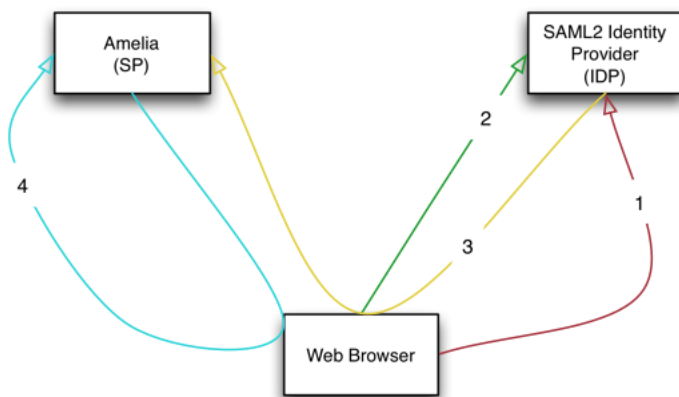


Figure 11. Amelia and Identity Provider Initiated Authentication

1. A user logs into the IDP directly. Normally by clicking a specially crafted link in some portal.

2. The user selects the service provider from a list provided by the IDP. This step will often be skipped if the user started with a special link in step 1.
3. The IDP generates and signed AuthnResponse and redirects the user's browser to Amelia.
4. Amelia verifies the signed response, and if valid, logs in the user and redirects the user's browser to the default page.

6.2 ANONYMOUS ACCESS

Anonymous access allows unauthenticated users to access Amelia. The unauthenticated user will be logged in as a special type of user, the Anonymous User. These users can be given access rights just like any other user in the system. The only difference between an anonymous user and any other user is that anonymous users will be auto-created when the user accesses the system and if anonymous access is enabled. Anonymous Users are given a default set of access rights based on a configured group. This group should have very minimal access rights. Most likely only End User for the anonymous domain.

Configuration of Anonymous access is a simple checkbox, which enables/disables the feature. This is also applied to the Amelia Escalation feature.

Using Anonymous Access are for specific use cases; please inquire with the Cognitive Team for specifics

6.3 SSL CERTIFICATES

Amelia requires SSL certificates to keep sensitive information sent across the network encrypted so that only the intended recipient can understand. IPsoft recommends wildcard certificates for ease of installation and future growth. If a wildcard SSL certificate is not used, four individual certificates or one "Subject Alternative Name", or SAN certificate must be created for each instance (except for DR).

Table 6. SSL Certificate Requirements

Name	Description
amelia.client.com	HTTP/HTTPS URL for Amelia

6.4 FIREWALL RULES

Table 7. Firewall Rules

Purpose	Source	Protocol	Destination	Destination Port	Required
Redirect to HTTPS	Browser	HTTP	haproxy	80	Yes
User interaction	Browser	HTTPS and WebSockets	haproxy	443	Yes
Statistics	Brower	HTTP	haproxy	1936	Yes
Antivirus	ipsoft-av-gateway	TCP	clamav	3310	Yes

Purpose	Source	Protocol	Destination	Destination Port	Required
Web UI	haproxy	HTTP (TLS 1.2)	amelia-engine-service@p001	4431	Yes
Monitoring	monitoring network	JMX	amelia-engine-service@p001	4434	Yes
RPC	haproxy	GRPC (TLS 1.2)	amelia-engine-service@p001	4435	Yes
RPC - All pods roundrobin	amelia-user-web	GRPC (TLS 1.2)	haproxy	4436	Yes
RPC - Engine pod 001 roundrobin.	amelia-user-web	GRPC (TLS 1.2)	haproxy	44001	Yes
Web UI	haproxy	HTTP (TLS 1.2)	amelia-user-web	4441	Yes
Monitoring	monitoring network	JMX	amelia-user-web	4444	Yes
UI Development	developer machine	HTTPS	webpack-dev-server	4449	For UI development only
Monitoring	monitoring network	JMX	amelia-escalation-service	4574	Yes
Monitoring	monitoring network	JMX	amelia-batch-service	4584	Yes
Web UI	haproxy	HTTP (TLS 1.2)	amelia-admin-web	4601	Yes
Monitoring	monitoring network	JMX	amelia-admin-web	4604	Yes
UI Development	deveveloper machine	HTTPS	webpack-dev-server	4609	For UI development only
Web UI	haproxy	HTTP (TLS 1.2)	amelia-admin-service	4611	Yes
Monitoring	monitoring network	JMX	amelia-admin-service	4614	Yes
Web UI	haproxy	HTTP (TLS 1.2)	amelia-engine-service@p002	4621	Yes
Monitoring	monitoring network	JMX	amelia-engine-service@p002	4624	Yes
RPC	haproxy	GRPC (TLS 1.2)	amelia-engine-service@p002	4625	Yes
RPC	amelia-user-web	GRPC (TLS 1.2)	haproxy	44002	Yes
Monitoring	monitoring network	JMX	amelia-integration-service	4634	Yes

Purpose	Source	Protocol	Destination	Destination Port	Required
RPC	haproxy	GRPC	amelia-integration-service	4635	Yes
Web UI	haproxy	HTTP (TLS 1.2)	amelia-account-service	4641	Yes
Monitoring	monitoring network	JMX	amelia-account-service	4644	Yes
RPC	haproxy	GRPC (TLS 1.2)	amelia-account-service	4645	Yes
RPC	amelia-*	GRPC (TLS 1.2)	haproxy	4646	Yes
Web UI	haproxy	HTTP (TLS 1.2)	amelia-model-server	4651	Yes
Monitoring	monitoring network	JMX	amelia-model-server	4654	Yes
RPC	haproxy	GRPC (TLS 1.2)	amelia-model-server	4655	Yes
RPC	amelia-*	GRPC (TLS 1.2)	haproxy	4656	Yes
REST Endpoints	haproxy	HTTP (TLS 1.2)	amelia-rest-gateway	4661	Yes
Monitoring	monitoring network	JMX	amelia-rest-gateway	4664	Yes
Monitoring	monitoring network	JMX	amelia-client-gateway	4674	Yes
Monitoring	monitoring network	JMX	amelia-client2-gateway	4684	Yes
Monitoring	monitoring network	JMX	amelia-client3-gateway	4694	Yes
Cluster Check	haproxy	HTTP	xinetd	13304	Yes
SQL	haproxy	TCP	mysql@amelia-kdb-master	13305	Yes
SQL	amelia-admin-service	TCP	haproxy	13306	Yes
SST	mysql@amelia-kdb-master	TCP	mysql@amelia-kdb-master	13307	Yes
Group Communication	mysql@amelia-kdb-master	TCP	mysql@amelia-kdb-master	13308	Yes
IST	mysql@amelia-kdb-master	TCP	mysql@amelia-kdb-master	13309	Yes
Slave Check	haproxy	HTTP	xinetd	13314	Yes
SQL	haproxy	TCP	amelia-kdb-slave	13315	Yes
SQL	amelia-engine-service	TCP	haproxy	13316	Yes
Cluster Check	haproxy	HTTP	xinetd	13324	Yes

Purpose	Source	Protocol	Destination	Destination Port	Required
SQL	haproxy	TCP	mysql@amelia-cdb-p001	13325	Yes
SQL	amelia-engine-service@p001	TCP	haproxy	13326	Yes
SST	mysql@amelia-cdb-p001	TCP	mysql@amelia-cdb-p001	13327	Yes
Group Communication	mysql@amelia-cdb-p001	TCP	mysql@amelia-cdb-p001	13328	Yes
IST	mysql@amelia-cdb-p001	TCP	mysql@amelia-cdb-p001	13329	Yes
Cluster Check	haproxy	HTTP	xinetd	13334	Yes
SQL	haproxy	TCP	mysql@amelia-cdb-p002	13335	Yes
SQL	amelia-engine-service@p002	TCP	haproxy	13336	Yes
SST	mysql@amelia-cdb-p002	TCP	mysql@amelia-cdb-p002	13337	Yes
Group Communication	mysql@amelia-cdb-p002	TCP	mysql@amelia-cdb-p002	13338	Yes
IST	mysql@amelia-cdb-p002	TCP	mysql@amelia-cdb-p002	13339	Yes
Cluster Check	haproxy	HTTP	xinetd	13344	Yes
SQL	haproxy	TCP	mysql@amelia-adb	13345	Yes
SQL	amelia-admin-service	TCP	haproxy	13346	Yes
SST	mysql@amelia-adb	TCP	mysql@amelia-adb	13347	Yes
Group Communication	mysql@amelia-adb	TCP	mysql@amelia-adb	13348	Yes
IST	mysql@amelia-adb	TCP	mysql@amelia-adb	13349	Yes
Datastore, Cache	amelia-user-web, amelia-admin-web, amelia-escalation	TCP	redis	13341	Yes
Redis Sentinel Cluster Management	redis-sentinel, haproxy	TCP	redis-sentinel	13342	Yes
Redis Sentinel VIP (watches all other redis clusters)	amelia-user-web, amelia-admin-web, amelia-escalation, amelia-engine, amelia-model-server	TCP	haproxy	13343	Yes
Messaging	haproxy	AMQP (TLS 1.2)	rabbitmq	13351	Yes

Purpose	Source	Protocol	Destination	Destination Port	Required
Messaging	amelia-*	AMQP (TLS 1.2)	haproxy	13352	Yes
Messaging	haproxy	STOMP (TLS 1.2)	rabbitmq	13353	Yes
Messaging	amelia-*	STOMP (TLS 1.2)	haproxy	13354	Yes
Monitoring	monitoring network	HTTP (TLS 1.2)	rabbitmq	13355	Yes
Clustering - Distribution	RabbitMQ	TCP	RabbitMQ	13356	Yes
Clustering - epmd	RabbitMQ/epmd	TCP	RabbitMQ/epmd	13357	Yes
Datastore, Cache	amelia-engine, amelia-model-server	TCP	redis	13361	Yes
Datastore, Cache	amelia-engine, amelia-model-server	TCP	redis	13371	Yes
Web UI	haproxy	HTTP (TLS 1.2)	amelia-model-server	13381	Yes
Monitoring	monitoring network	JMX	amelia-model-server	13384	Yes
RPC	haproxy	GRPC (TLS 1.2)	amelia-model-server	13385	Yes
RPC	amelia-*	GRPC (TLS 1.2)	haproxy	13386	Yes
Date-time parser	amelia-engine-service	HTTPS	haproxy	14010	Yes
Date-time parser	haproxy	HTTPS	amelia-duckling-service	14011	Yes
AV scan REST API	amelia-*-web	HTTP (TLS 1.2)	haproxy	14020	Yes
AV scan REST API	haproxy	HTTP (TLS 1.2)	ipsoft-av-gateway	14021	Yes
Monitoring	monitoring network	JMX	ipsoft-av-gateway	14024	Yes
syntaxnet parser	haproxy	TCP	amelia-tf-syntaxnet	14000	Yes
syntaxnet parser	amelia-*	GRPC	amelia-tf-syntaxnet	14001	Yes
syntaxnet parser	amelia-tf-syntaxnet	TCP	amelia-tf-syntaxnet	14009	Yes
syntaxnet parser	haproxy-test	TCP	amelia-tf-syntaxnet-test	14090	Yes
syntaxnet parser	amelia-*-test	GRPC	amelia-tf-syntaxnet-test	14091	Yes
syntaxnet parser	amelia-tf-syntaxnet-test	TCP	amelia-tf-syntaxnet-test	14099	Yes
Flow UI	haproxy	HTTP	amelia-h2o	54321	Yes

Purpose	Source	Protocol	Destination	Destination Port	Required
H2O Internal Communication	amelia-h2o	TCP	amelia-h2o	54322	Yes

7. Monitoring

For IPsoft Hosted and Online deployments, all Amelia instances and supporting Operating Systems are monitored by IPcenter using IPmons. IPmons are configured to monitor each component of Amelia; each component having several individual checks to ensure thorough reporting and availability of each Amelia instance. Below is an example of the OS checks deployed:

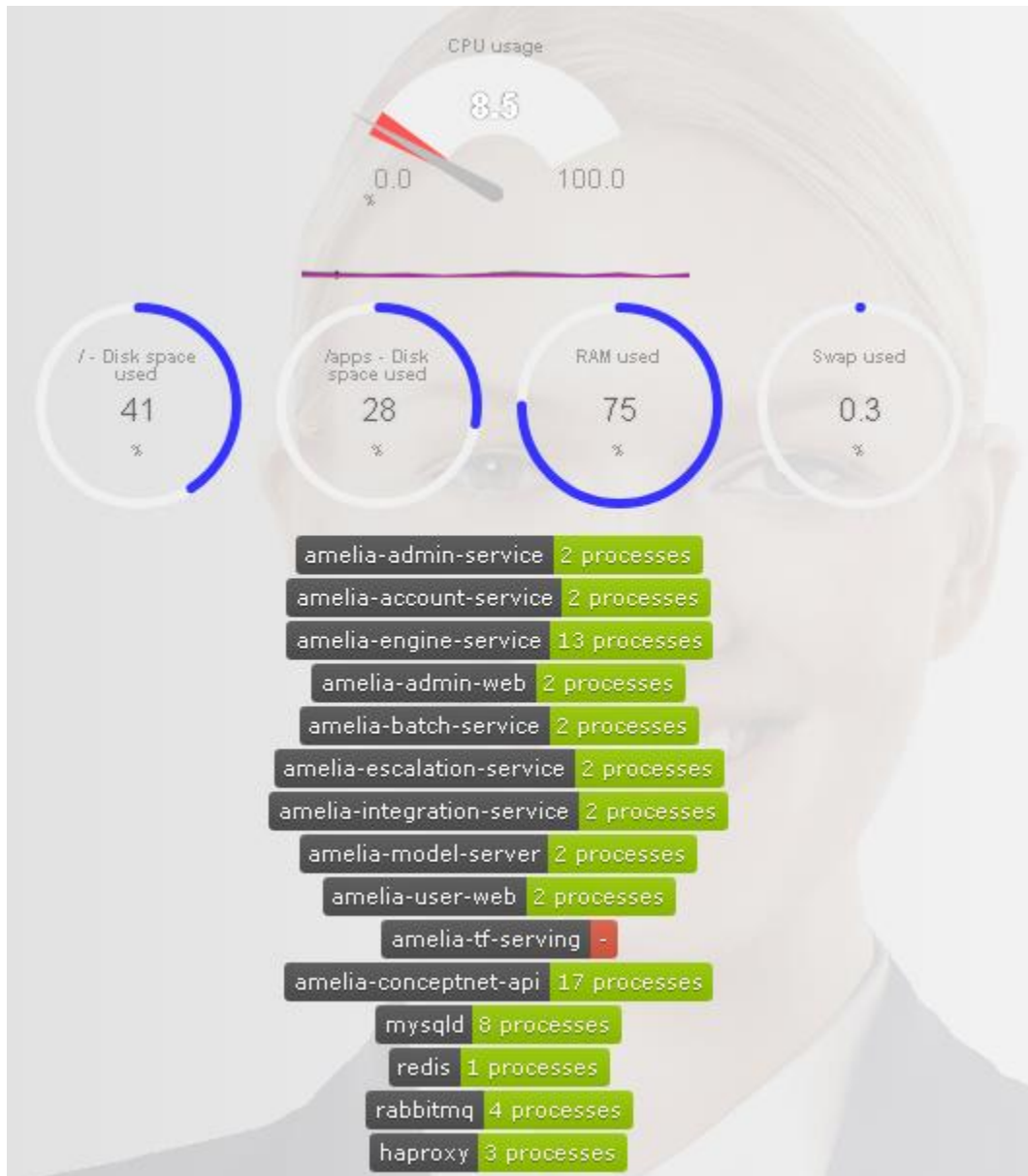
Service Status Details for Host 'app01.dev.amelia.ipcenter.com.ipsoft'

Limit Results: 100

Host	Service	Status	Last Check	Duration	Attempt	Status Information
app01.dev.amelia.ipcenter.com.ipsoft	Amelia Backup Log - /apps/backup	OK	02-07-2017 15:33:41	5d 1h 29m 29s	1/1	OK: Log file is clean.
	Amelia SSL Expiry	OK	02-07-2017 15:33:57	40d 23h 32m 41s	1/3	OK: SSL certificate on app01.dev.amelia.ipcenter.com:443 expires in 943.14 days. Wa
	Amelia Webcheck	OK	02-07-2017 15:32:27	55d 5h 39m 11s	1/1	OK: Total time: 0.1234 - All Checks Passed
	Amelia YUM Audit Log	OK	02-07-2017 15:33:41	20d 4h 7m 29s	1/1	OK: Log file is clean.
	Disk - /	OK	02-07-2017 15:32:07	99d 4h 53m 43s	1/3	OK: Disk usage for "/": 4.54GB (22%); used: 4.54GB, free: 15.45GB, total: 19.99GB, var
	Disk - /apps	OK	02-07-2017 15:32:07	99d 4h 53m 43s	1/3	OK: Disk usage for "/apps": 11.71GB (23%); used: 11.71GB, free: 38.26GB, total: 49.97
	Disk - /boot	OK	02-07-2017 15:32:07	99d 4h 53m 43s	1/3	OK: Disk usage for "/boot": 175.57MB (35%); used: 175.57MB, free: 321.09MB, total: 497.67MB
	Disk - /dev/shm	OK	02-07-2017 15:32:07	99d 4h 53m 43s	1/3	OK: Disk usage for "/dev/shm": 12.00KB (0%); used: 12.00KB, free: 31.38GB, total: 31.39GB
	Disk - /run	OK	02-07-2017 15:32:07	99d 4h 53m 43s	1/3	OK: Disk usage for "/run": 600.41MB (1%); used: 600.41MB, free: 30.79GB, total: 31.39GB
	Disk - /sys/fs/cgroup	OK	02-07-2017 15:32:07	99d 4h 53m 43s	1/3	OK: Disk usage for "/sys/fs/cgroup": 0B (0%); used: 0B, free: 31.38GB, total: 31.39GB
	Disk - /tmp	OK	02-07-2017 15:32:07	20d 4h 15m 19s	1/3	OK: Disk usage for "/tmp": 16.00KB (0%); used: 16.00KB, free: 31.38GB, total: 31.39GB
	Disk - /var/tmp	OK	02-07-2017 15:32:07	20d 4h 15m 19s	1/3	OK: Disk usage for "/var/tmp": 16.00KB (0%); used: 16.00KB, free: 31.38GB, total: 31.39GB
	Disk Latency	OK	02-07-2017 15:34:34	0d 3h 4m 13s	1/3	OK: dm-0 await: 0.00 dm-1 await: 0.00 dm-2 await: 0.00 sda await: 0.00 sdb await: 0.00
	HAProxy Back End	OK	02-07-2017 15:30:25	50d 3h 9m 45s	1/3	OK: Connection utilization for stats: 0.00, thresholds 50/80, stats status: UP, expected U
	HAProxy Front End	OK	02-07-2017 15:30:25	50d 3h 9m 45s	1/3	OK: Connection utilization for stats: 0.00, thresholds 50/80, stats status: OPEN, expecte
	Host Memory	OK	02-07-2017 15:32:25	99d 4h 50m 24s	1/3	OK: Memory utilization: 35.35%. Warning/Critical thresholds: 95/98
	IPremoted	OK	02-07-2017 15:31:57	99d 4h 53m 1s	1/3	OK: 0.0080 sec. response time. "IPremote" matched, received "IPremote - 5.5.3"
	Inodes - /	OK	02-07-2017 15:32:00	99d 4h 50m 45s	1/3	OK: Disk inode usage for "/": 82677 (1%); used: 82677, free: 20888843, total: 2097152
	Inodes - /apps	OK	02-07-2017 15:32:00	99d 4h 50m 45s	1/3	OK: Disk inode usage for "/apps": 5833 (1%); used: 5833, free: 52418871, total: 52424
	Inodes - /boot	OK	02-07-2017 15:32:00	99d 4h 50m 45s	1/3	OK: Disk inode usage for "/boot": 340 (1%); used: 340, free: 511660, total: 512000, wa
	Inodes - /dev/shm	OK	02-07-2017 15:32:00	99d 4h 50m 45s	1/3	OK: Disk inode usage for "/dev/shm": 4 (1%); used: 4, free: 8225356, total: 8225360, v
	Inodes - /run	OK	02-07-2017 15:32:00	99d 4h 50m 45s	1/3	OK: Disk inode usage for "/run": 463 (1%); used: 463, free: 8224897, total: 8225360, v
	Inodes - /sys/fs/cgroup	OK	02-07-2017 15:32:00	99d 4h 50m 45s	1/3	OK: Disk inode usage for "/sys/fs/cgroup": 13 (1%); used: 13, free: 8225347, total: 822
	Linux Messages Log	OK	02-07-2017 15:33:41	104d 23h 23m 11s	1/1	OK: Log file is clean.
	Load Average	OK	02-07-2017 15:34:07	99d 4h 51m 26s	1/5	OK: load average: 0.01,0.03,0.05; wload: 9999.00,24.00,9999.00 dload: 9999.00,32.00,
	Perf Data	OK	02-07-2017 15:32:25	99d 4h 50m 24s	1/3	OK: All perfdata retrieved
	Proc - crond	OK	02-07-2017 15:33:44	99d 4h 54m 2s	1/3	OK: 1 process running with arguments: "/usr/sbin/crond -n", as regex
	Proc - haproxy	OK	02-07-2017 15:33:44	49d 22h 56m 20s	1/3	OK: 2 processes running with arguments: "/usr/sbin/haproxy -f /etc/haproxy/haproxy.cfg

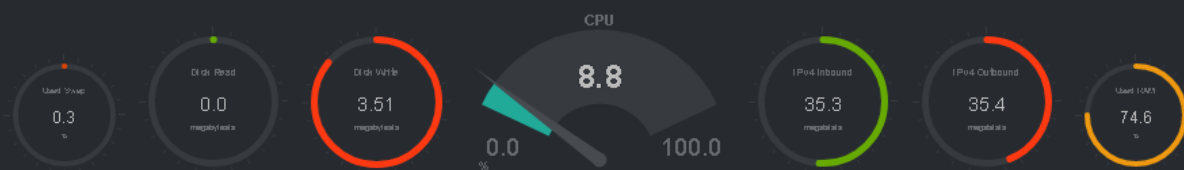
Figure 12. Example of OS Monitoring Checks Performed

Starting with Amelia v3, IPsoft is distributing *netdata* for insights into system and application performance. It also provides a rudimentary level of monitoring and alerting. These alerts do not currently come back to IPsoft, but may be directed to a client's email if desired.



System Overview

Overview of the key system metrics.



8. Backups

Backups for Amelia can be outfitted with the Customer/Partner's Enterprise Backup solution. This allows the customer to roll Amelia backups into their current enterprise policies. Although most deployments leverage a virtual environment, backups of entire VM can cause a degraded or unresponsive system; it is highly recommended to disable the "virtual machine memory" and "Quiesce guest file system".

File level backups would require a more unique solution. This would require standing up a new, fresh install of an Amelia host to replace an unrecoverable one, restoring the original's configuration and data from a remote backup. For a complete list of backup locations and uses, please contact the Service Technology team.

Defining Recovery Time Objective (RTO), Recovery Point Objective (RPO), and Geographic Redundancy Objective (GRO) should be considered when creating backup/recovery policy for Amelia. For IPSoft Hosted Instances, non-production instances have a 24-hour RTO and no larger than 24-hour RPO. For production instances, a 4-hr RTO and a 4-hour RPO.

9. Disaster Recovery

In its current architecture, Amelia v3 current supports an Active/Passive approach, where web traffic would be directed to the "live" datacenter and replicated to a secondary datacenter. The replication method can be achieved at either the middleware (native) or infrastructure level.

For native replication, Percona XtraDB Cluster would be replicating the databases across the WAN to the DR databases. In this setup all Production and DR database nodes are configured as one large logical cluster. Newer versions of Amelia, Language Packs, and Gateways would be performed separately for both the Production and DR instances; not as one upgrade to cover both instances.

Infrastructure replication can be achieved using 3rd party hardware/software vendors. SAN based replication with orchestration tools (such as Zerto/VMware SRM) is a proven DR solution, as well as software based solutions such as Veeam and Veritas. A Production instance is replicated without changing the hostnames of the VMs, however it is best to keep the IP addresses identical if possible using a stretched layer 2 network.

The length of time necessary to conduct a failover of Amelia will depend on the overall design and available automation processes. Regarding Amelia specifically, all Amelia processes can be started in parallel and can take up to 5 minutes for Amelia to be started in the secondary datacenter, assuming IP addresses are not altered.

For IPSoft Hosted instances, a 4-hr RTO and a 4-hour RPO.

10.High Availability Configurations

A single Amelia deployment has can be scaled up to 8 conversation PODs, each handling 150 concurrent users for a total for 1200 concurrent conversations.

Each conversation POD contains 3 Amelia Conversation Engines and related services including Redis for session data, a Percona CDB database cluster, and a KDB Database slave.

Additional services such as user and Admin web front ends and Admin services are deployed on additional virtual machines.

This configuration has been deployed in production handling >200,000 daily conversations with end users.

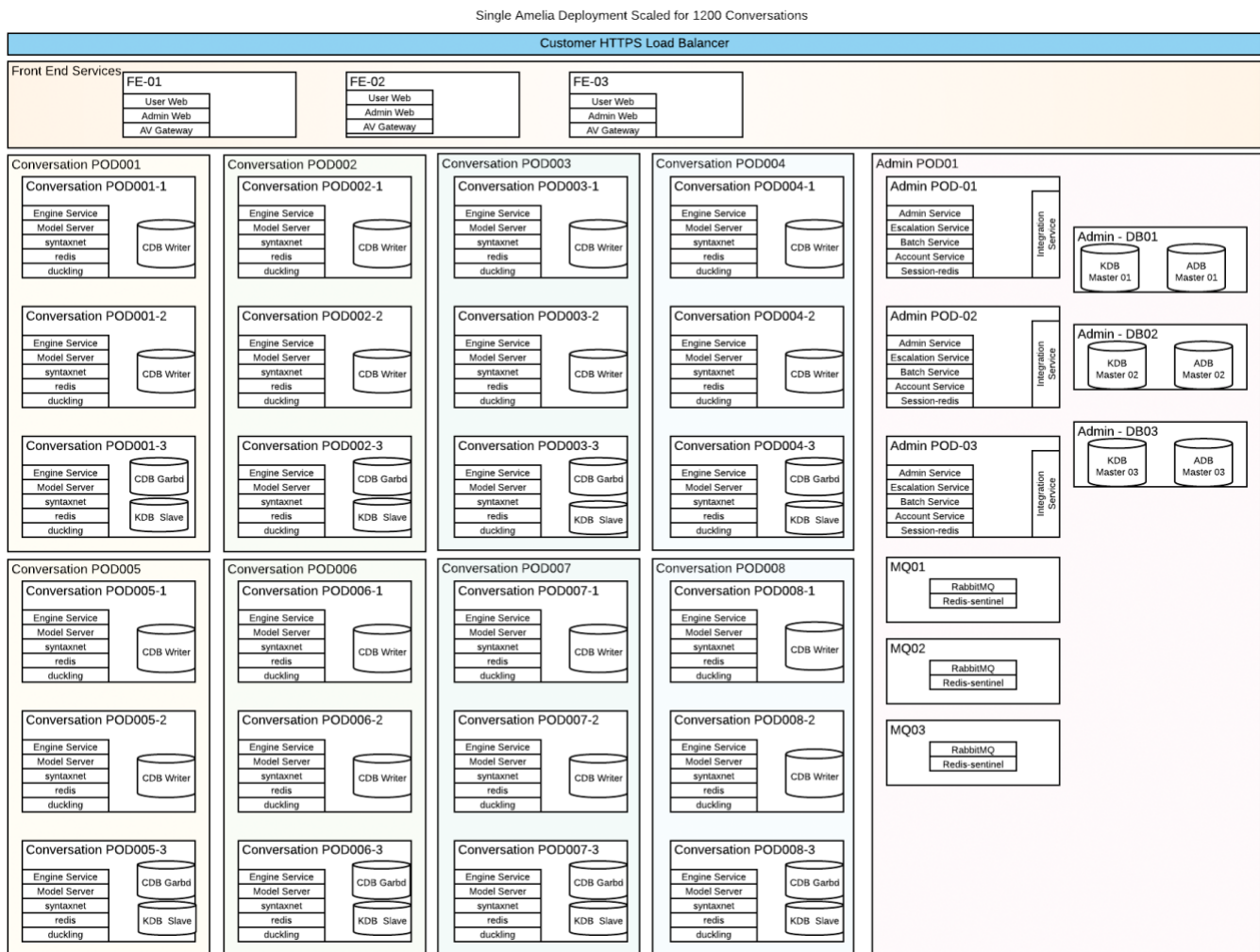


Figure 13. Amelia Deployment Scaled for 1200 Concurrent Conversations

10.1 FRONT END SERVICES

The front end services include the User Web Services, Admin Web Services and AntiVirus Gateway

Each of the front end service hosts are configure with a minimum of:

CPU	Memory	Disk
8 vCPU	32 GB Ram	OS + 100GB for application and logs

10.2 ADMIN POD

The Admin POD contains all admin services such a escalation, integration, batch and the account services as well as dedicated hosts for the KDB and ADB Percona clusters and dedicated RabbitMQ Hosts.

The Rabbit hosts also run Redis Sentinel which is responsible for Redis cluster management throughout the deployment.

Each of the Admin hosts are configured with a minimum of:

CPU	Memory	Disk
8 vCPU	64 GB Ram	OS + 300GB for application and logs

Each of the Admin Database hosts are configured with a minimum of:

CPU	Memory	Disk
8 vCPU	64 GB Ram	OS + 1TB for databases, backups and logs

10.3 RABBITMQ HOSTS

Each of the Admin RabbitMQ hosts are configured with a minimum of:

CPU	Memory	Disk
8 vCPU	16 GB Ram	OS + 100GB for applications and logs

10.4 CONVERSATION POD

Each node in the conversation pod contains the Amelia engine as well as associated components such as the Model Server, Syntaxnet, Redis (for session data) and the duckling service.

The first two nodes contain 2/3 of the CDB cluster with 2 writers (only one is actively used at a time). The 3rd node contains a CDB Garbd for maintaining cluster integrity as well as a KDB slave used for reading knowledge during conversations.

Each of the Conversation POD hosts are configured with a minimum of:

CPU	Memory	Disk
8 vCPU	64 GB Ram	OS + 300GB applications, databases, backups and logs

10.5 GENERAL NOTES ABOUT HARDWARE/RESOURCE RECOMMENDATIONS

10.5.1 CPU

- All Virtual CPUs are mapped to exactly 1 CPU core of an Intel E5-2687W v3 @ 3.10GHz or greater CPU.
- No over-subscription is permitted. CPUs must have a base frequency of 3Ghz.

10.5.2 Disk

- Disk space can be influenced by the number of conversations that are handled and size of each conversation and retention times. Observe average growth per day and plan accordingly.
- It is suggested to utilize LVM for all application file systems to allow for easy growth as required.
- Flash/SSD/NVMe based storage is required for Amelia to operate effectively.
- IO Subsystems must be able to sustain a minimum of 1000 Mixed IOPS (80% write, 20% Read) per VM under load with a <10ms response time to read and write requests.
- Each VM will perform an average of 50MBps of write activity with an additional 10MBps of read, when Journey Analytics and detailed history are disabled. IPsoft recommends sizing for 100MBps of Writes and 25MBps of reads per VM.
- When Journey Analytics and detailed history are enabled, it is suggested to triple the disk IO figures.

10.5.3 VMware

- VMware's vMotion has the potential to adversely affect Amelia's performance and stability in certain situations. This includes synchronization of the RabbitMQ cluster as well as Percona database clustering.
- It is suggested that if vMotion is required, it be scheduled to execute during quiet hours of operation.
- DRS can have the same effect as it used vMotion underneath. It is suggested to disable fully automated DRS or set to a less aggressive level.

10.5.4 OS

- Each VM has a minimum of 16GB Swap

10.6 SCALING AMELIA COMPONENTS

Amelia was designed to be able to run each of its components on different OS Images for scaling/distribution of load and to meet security requirements.

10.6.1 Databases

Databases may be isolated to dedicated hosts if required by information security requirements. This includes KDB, ADB and CDB.

10.6.2 Scaling ADB and KDB

The reference architecture describers here is scaled appropriately for the 1200 concurrent user base. No further scaling should be necessary.

10.6.3 Scaling CDB

We believe the proper distribution of services and databases for the CDB requirement is captured in the reference design.

All databases are heavily used for both reads and writes (hence the minimum of 1000 IOPS requirement listed above) during conversations and usage depends on the number of conversations occurring in the given POD

10.6.4 Alternate CDB Design (Not recommended)

IPsoft has tested performance of Amelia using 3 conversations PODs and 3 dedicated nodes for the CDB Database cluster in a 3x3 configuration

Node 1	Node 2	Node 3
CDB01-Master	CDB01-Writer	CDB01-GarbD
CDB02-Writer	CDB02-GardB	CDB02-Master
CDB03-GarbD	CDB03-Master	CDB03-Writer

This design appears to offer a limited to no benefit of lower CPU utilization on the conversation POD nodes. However, the CDB nodes usage becomes very CPU and network intensive due to the number of conversations being processed.

In a shared VMware environment the overall utilization remained constant in both configurations under load.

Currently, testing reveals no benefit from a conversations per POD perspective, or a conversation mean time to response perspective. Therefore, we do not recommend this configuration.

10.6.5 Rabbit MQ

Dedicated RabbitMQ/Sentinel nodes are provided in the reference architecture to provide for necessary resources response time required by these services.

10.6.6 Integration Services

Integration services run with the rest of the admin services. We have observed up to 30 integration flows per second run in the configuration with a 75ms response time (inclusive of the back end systems).

In a 1000 concurrent user configuration with a heavy reliance on integrations we generally see an average of 10 flows per second executing with minimal load on the integration servers.

Integration services may be scaled to their own VMs and may be added in batches of 3 nodes to provide additional scaling.

Scaling of integration services is beyond the scope of this testing due to the variety of back end systems that may influence integration performance.