

Kadeploy3

Efficient and Scalable Operating System Provisioning for Clusters

EMMANUEL JEANVOINE, LUC SARZYNIEC, AND LUCAS NUSSBAUM



Emmanuel Jeanvoine is a Research Engineer at Inria Nancy Grand Est. He specializes in distributed systems and high performance computing. In 2007, he obtained a PhD in computer sciences at Université de Rennes 1. emmanuel.jeanvoine@inria.fr



Luc Sarzyniec is a Junior Engineer at Inria Nancy Grand Est. He started working on distributed systems and high performance computing after he obtained his master of research in computing sciences in 2011. luc.sarzyniec@inria.fr



Lucas Nussbaum is an Assistant Professor at Université de Lorraine. His research focuses on high performance computing and distributed systems. He is also involved in a professional curriculum focusing on system administration. lucas.nussbaum@loria.fr

Installing an operating system can be tedious when it must be reproduced on many computers, on large scale clusters, for instance. Because installing the nodes independently is not realistic, disk cloning or imaging with tools such as Clonezilla [1], Rocks [5], SystemImager [6], or xCAT [8] is a common approach. In those cases, the administrator must keep updated just one node (sometimes called the golden node) that will be replicated to other nodes. In this article, we present Kadeploy3, a tool designed to perform operating system provisioning using disk imaging and cloning. Thanks to its efficiency, scalability, and reliability, this tool is particularly suited for large scale clusters.

Reliable Deployment Process with Kadeploy3

Kadeploy3 belongs to the family of disk imaging and cloning tools. It takes as input an archive containing the operating system to deploy, called an environment, and copies it on the target nodes. As a consequence, Kadeploy3 does not install an operating system following the classical installation procedure, and the user must provide an archive of the environment (as a tarball, for Linux environments).

Kadeploy3 does not directly take control of the nodes because doing so requires some specific and uncommon hardware support. Instead, it uses common network boot capabilities based on the PXE protocol [4], and it manages the associated PXE profiles.

Using such a mechanism, combined with the capability to update the PXE profiles of the nodes dynamically and to reboot the nodes in a reliable way (thanks to out-of-band control interfaces, such as Baseboard Management Controller, Remote Supervisor Adapter, or the power distribution unit's capabilities), taking control of the nodes and specifying what they are booting is possible.

As shown in Figure 1, a typical deployment with Kadeploy3 is composed of three major steps, called macro steps.

1. Minimal environment setup: the nodes reboot into a trusted minimal environment that contains all the tools required for the deployment (partitioning tools, archive management, etc.), and the required partitioning is performed.
2. Environment installation: the environment is broadcast to all nodes and extracted on the disks. Some post-installation operations also can be performed.
3. Reboot using the newly deployed environment.

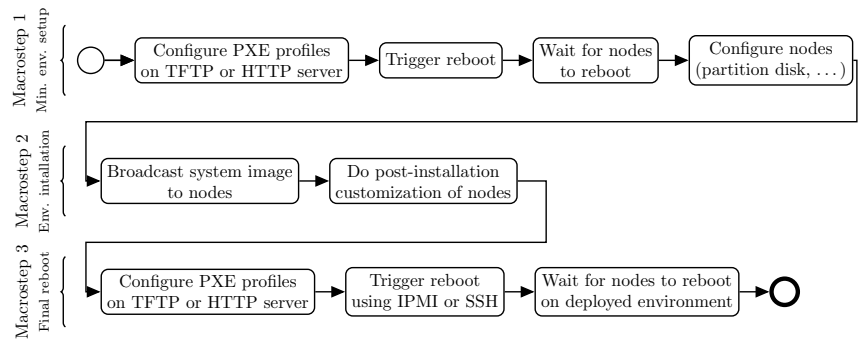


Figure 1: Kadeploy deployment process, composed of three macro-steps

Each macro step can be executed via several different mechanisms to optimize the deployment process depending on required parameters and the specific infrastructure. For instance, the reboot using the newly deployed environment step can perform a traditional reboot or it might instead rely on a call to `kexec(8)` for a shorter reboot.

Reconfiguring a set of nodes involves several low-level operations that can lead to failures for various reasons, e.g., temporary loss of network connectivity, reboot taking longer than planned, etc. Kadeploy3 reliability is achieved because (1) the deployment process has powerful error management and (2) critical reboot operations required for the node control are based on reboot commands escalation in order to be able to take control of the nodes in any situation.

Reliability of the Deployment

Kadeploy3 is designed to detect failures as quickly as possible and improve deployment reliability by providing a macro-step replay mechanism on the nodes of interest. To illustrate that, let's consider the last deployment macro step that aims at rebooting using the deployed environment. Kadeploy3 implements, among others, the following strategies:

1. Directly load the kernel inside the deployed environment thanks to `kexec`.
2. Perform a hard reboot using out-of-band management hardware without checking the state of the node.

Thus it is possible to describe strategies such as: try the first strategy; if some nodes fail, try the second strategy, several times if required.

Because all the steps involved in the deployment process rely on system calls (hard disk operations, network communications, specific hardware management), special attention has been paid to error handling. Kadeploy3 collects the result of every operation (exit status, stdout, stderr), even when it is performed on remote nodes. As a consequence, some steps can be replayed on nodes where a problem occurs.

Furthermore, some operations may last too long (e.g., network boot, file-system creation, etc.), but Kadeploy3 provides administrators with the capability of defining specific timeouts for some operations in order to adapt the deployment process to the infrastructure. That allows identifying some problems quickly and replaying some operations on the related nodes.

Reliability of Reboot Operations

Because reboot operations are essential to control the cluster nodes, and ultimately the entire deployment process itself, they must behave correctly and reliably. Several methods can be used to reboot nodes, for instance:

1. Directly execute the `/sbin/reboot` command.
2. Use out-of-band management hardware with protocols such as IPMI. Various kinds of reboots can be executed: reset, power cycle, etc.
3. Use the power management capability of the power distribution unit (PDU).

Performing an `/sbin/reboot` is the best solution with regards to speed and cleanliness; however, it may not be an option if the target node is unreachable via in-band methods such as SSH (e.g., the node is already down, the OS has crashed, an unfriendly operating system is installed, etc.). In this scenario, we would use IPMI-like features if available. Also, because it bypasses the power-on self test, it might be better for speed to perform a reset rather than a power cycle, but sometimes this is not sufficient. Finally, if onboard management hardware is unreachable, we may be required to use the capabilities of a remotely manageable PDU.

Kadeploy3 provides administrators with a way to specify several levels of commands in order to perform escalation if required. This allows them to perform highly reliable deployments if the clusters have the appropriate hardware. Unfortunately, depending on the methods chosen, reboot escalation comes at a cost, and a balance must be struck between desired reliability and the time to deployment.

Scalability

In addition to having a reliable node-control mechanism, deploying large scale clusters in a reasonable time requires being able to execute several commands efficiently and to send large files on a large number of nodes.

Parallel Commands

The deployment workflow contains several operations that reduce to executing a command on a large set of nodes.

Thanks to SSH, one can execute commands remotely and retrieve their outputs, but launching SSH commands on a large number of nodes in sequence does not scale at all. Furthermore, launching all commands simultaneously can impose an extreme load on the server and can consume all of its file descriptors.

Several tools have been built to overcome these limitations. For instance, Pdsh [3] and ClusterShell [2] are designed to execute SSH commands on many nodes in parallel. Both tools use windowed execution to limit the number of concurrent SSH commands, and both also allow retrieval of command outputs on each node.

We choose to leverage TakTuk [9] as our mechanism for parallel command execution and reporting. TakTuk is based on a model of hierarchical connection. This allows TakTuk to distribute the execution load on all the nodes in a tree and to perform commands with low latency. Using such a hierarchical mechanism would normally require the tool to be installed on all nodes. Fortunately, TakTuk includes a convenient auto-propagation feature that ensures the tool's existence on all necessary nodes. The tool also uses an adaptive work-stealing algorithm to improve performance, even on heterogeneous infrastructures.

File Broadcast

The broadcast of the system image to all nodes is a critical part of the deployment. In cluster environments where the most important network for applications is using Infiniband or Myrinet, the Ethernet network is often composed of a hierarchy of switches (e.g., one switch per rack) that is hard to leverage for a high-performance broadcast. File distribution to a large number of nodes via any sequential push or pull method is not scalable. Kadeploy3 provides system administrators with three scalable file distribution approaches during the *Environment installation* macro step to minimize deployment time.

With tree-based broadcast, a file is sent from the server to a subset of nodes, which in turn send the file to other subsets until all the nodes have received the file. The size of the subsets, called tree arity, can be specified in the configuration. A large arity can reduce the latency to reach all nodes, but transfer times might increase because global bandwidth is equal to the bandwidth of a network link divided by the tree arity. The opposite effect occurs when the arity is small. In general, this broadcast method does not maximize bandwidth and should be used primarily for the distribution of small files. This method is also inefficient when used in hierarchical networks. We implement tree-based broadcast using TakTuk.

Chain-based broadcast facilitates the transfer of files with high bandwidth. A classical chain-based broadcast suffers from the establishment time of the chain in large-scale clusters. Indeed, because each node must connect to the next node in the chain (usually via SSH), a sequential initialization would drastically increase the entire broadcast period. Thus we perform the initialization of the chain with a tree-based parallel command. This kind of broadcast is near-optimal in a hierarchical network if the chain is well ordered because, as shown in Figure 2, all the full-duplex network links can be saturated in both directions, and the performance bottleneck becomes the backplane bandwidth of the network switches. For this method, we implement chain initialization using TakTuk and perform transfers using other custom mechanisms.

BitTorrent-based broadcast is able to send files at large scale without making any assumptions about the quality of the network. Furthermore, BitTorrent is able to handle churn efficiently, an important property in large scale systems such as petascale and future exascale clusters. Currently, our experiments show that there are two scenarios in which the performance of this broadcast method is inferior to the other methods. The first pathological case is one in which we are broadcasting

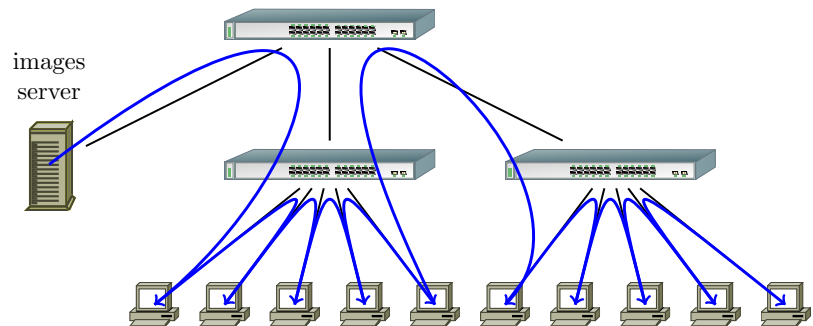


Figure 2: Topology-aware chained broadcast. Data is pipelined between all nodes. When correctly ordered, this ensures that inter-switch links are only used once in both directions.

on a small-scale cluster with a high-speed network, and the second is one in which we are broadcasting small files. In both cases, BitTorrent exhibits high latency, and the overhead of the protocol dominates the time to broadcast. The large number of established connections between nodes induced by the protocol can lead to bottlenecks depending on the network topology.

In a default configuration, Kadeploy3 uses tree-based broadcast for the files used in the deployment process (e.g., disk partition map) and chain-based method for the environment broadcast that is usually a large file; however, this behavior can be modified in the configuration.

Other Advanced Features

In addition to being reliable and scalable, Kadeploy has many useful features.

Multi-Cluster Support

Kadeploy3 can be configured to manage several clusters at the same time through a hierarchical set of YAML configuration files. In a grid-like environment, initiating and controlling deployments on several Kadeploy servers from a unique Kadeploy client is also possible.

Hardware and Software Compatibility

Kadeploy3 does not generally rely on vendor-specific mechanisms. Vendor-specific remote control systems used to trigger node reboots can be used easily, even if they do not support the IPMI protocol. Environments can be stored either as tarballs (for Linux environments) or as raw partitions, which enables the deployment of Windows or BSD-based systems.

Rights Management and Environments Library

Kadeploy3 can be used to provide users with a cloud-like experience with bare-metal system reservation. It can integrate with a cluster batch scheduler used to manage reservations in order to delegate system deployment rights to specific users for the duration of a job. Kadeploy3 can also manage a set of environments and their visibility (public, private) in order to provide default environments, on which users can base their work to create and register custom environments.

Statistics Collection

Identifying defunct nodes in a cluster is often hard, especially when failures are transient. Kadeploy3 integrates a statistics-collection mechanism that enables the detection of nodes that often fail during the deployment process.

Performance Evaluation

Grid'5000 Experimental Testbed

Kadeploy3 has been used intensively on the Grid'5000 testbed (<http://www.grid5000.fr>) since the end of 2009 (and previous versions of Kadeploy were used since 2004). In that time, approximately 620 different users have performed 117,000 deployments. On average, each deployment has involved 10.3 nodes. The largest deployment involved 496 nodes. To our knowledge, the deployed operating systems are mostly based on Linux (all flavors) with a sprinkling of FreeBSD.

Although the Grid'5000 use case does not exercise all the goals targeted by Kadeploy3 (e.g., scalability), it shows the tool's adequacy with regard to most characteristics, such as reliability.

Curie Petascale Supercomputer

We had the opportunity to evaluate Kadeploy3 on the Curie [7] supercomputer owned by GENCI (<http://www.genci.fr/>) and operated by CEA (<http://www.cea.fr/>): 2088 nodes were available to perform the test and the goal was to deploy the production environment. After a single administrative cycle, 2015 nodes were successfully deployed. This proved the efficiency and the reliability of Kadeploy3 in a large-scale production infrastructure.

Virtual Testbed

Validating scalability on large physical infrastructures can become complex because it requires privileged rights on many components (e.g., access to management cards, modification of PXE profiles, etc.). For example, because the Curie supercomputer is used for production purposes, we only had access to it for several hours. Thus we chose to build our own large-scale virtual testbed on Grid'5000, leveraging important features such as link-layer isolation, and Kadeploy3 of course.

We performed an experiment in which we used 635 physical nodes of the Grid'5000 testbed. Depending on the nodes capabilities, we launched a variable number of KVM virtual machines. In total, 3,999 virtual machines were launched and participated in a single virtual network (despite that the physical nodes were located on four different sites). Then we installed all the required servers: DHCP, TFTP, MySQL, HTTP server, Kadeploy3. Once the testbed was launched, we were able to perform deployments within a single cluster of 3,999 nodes. During the largest run, a 430 MB environment was installed on 3,838 virtual machines in less than an hour; 161 virtual nodes were lost due to network or KVM issues. A significant amount of time was also wasted because of the high latency between geographically distant sites (10–20 ms), which affected some infrastructure services such as DHCP and the PXE protocol.

Wrapping Up

We think that Kadeploy3 can help system administrators of large-scale clusters save precious time by reducing OS provisioning time. The best way to be convinced is to try it. Kadeploy3 is free software (CeCill 2 license) written in Ruby and available from <http://kadeploy3.gforge.inria.fr/>. Source code, as well as Debian and RPM packages, can be downloaded. Kadeploy3 is configured thanks to few YAML files. To help administrators, a complete guide describes the entire installation and configuration process [10].

References

- [1] Clonezilla: <http://clonezilla.org>.
- [2] ClusterShell: <http://cea-hpc.github.com/clustershell>.
- [3] Parallel Distributed Shell: <http://sourceforge.net/projects/pdsh>.
- [4] Preboot Execution Environment (PXE) Specification: <http://download.intel.com/design/archives/wfm/downloads/pxespec.pdf>.

[5] Rocks: Open-Source toolkit for real and virtual clusters: <http://www.rocksclusters.org>.

[6] SystemImager: <http://systemimager.org>.

[7] The Curie supercomputer: <http://www-hpc.cea.fr/en/complexe/tgcc-curie.htm>.

[8] xCAT: Extreme Cloud Administration Toolkit: <http://xcat.sourceforge.net>.

[9] Benoit Claudel, Guillaume Huard, and Olivier Richard, "TakTuk: Adaptive Deployment of Remote Executions," Proceedings of the 18th ACM International Symposium on High Performance Distributed Computing (HPDC), ACM 2009, pp. 91-100.

[10] Kadeploy3: <https://gforge.inria.fr/frs/download.php/27606/kadeploy-3.1-3.pdf>.

