

# Searching for High-Fidelity Builds Using Active Learning

Harshitha Menon\*, Konstantinos Parasyris\*, Tom Scogland, Todd Gamblin

{harshitha,parasyris1,scogland1,tgamblin}@llnl.gov

Lawrence Livermore National Laboratory

Livermore, CA, USA

## ABSTRACT

Modern software is incredibly complex. A typical application may comprise hundreds or thousands of reusable components. Automated package managers can help to maintain a consistent set of dependency versions, but ultimately the solvers in these systems rely on constraints generated by humans. At scale, small errors add up, and it becomes increasingly difficult to find high-fidelity configurations. We cannot test all configurations, because the space is combinatorial, so exhaustive exploration is infeasible.

In this paper, we present *Reliabuild*, an auto-tuning framework that efficiently explores the build configuration space and learns which package versions are likely to result in a successful configuration. We implement two models in *Reliabuild* to rank the different configurations and use adaptive sampling to select good configurations with fewer samples. We demonstrate *Reliabuild*'s effectiveness by evaluating 31,186 build configurations of 61 packages from the Extreme-scale Scientific Software Stack (E4S). *Reliabuild* selects good configurations efficiently. For example, *Reliabuild* selects 3× the number of good configurations in comparison to random sampling for several packages including Abyss, Bolt, libnrm, OpenMPI. Our framework is also able to select all the high-fidelity builds in half the number of samples required by random sampling for packages such as Chai, OpenMPI, py-petsc4py, and slepc. We further use the model to learn statistics about the compatibility of different packages, which will enable package solvers to better select high-fidelity build configurations automatically.

## ACM Reference Format:

Harshitha Menon\*, Konstantinos Parasyris\*, Tom Scogland, Todd Gamblin. 2022. Searching for High-Fidelity Builds Using Active Learning. In *19th International Conference on Mining Software Repositories (MSR '22)*, May 23–24, 2022, Pittsburgh, PA, USA. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3524842.3528464>

## 1 INTRODUCTION

Since at least the 1960's, developers have striven to make use of software components [38]. Reusing components saves time and separates concerns—client code can rely on robust implementations of common functionality without reimplementing them. However, with the efficiency of reuse has come an increase in complexity [10, 16, 19, 26, 36]. Package management is a cornerstone of

modern software engineering; millions of components, or *packages* are available from public registries and can be included in a project by simply running a command or modifying a line in a file.

To deal with this complexity, software ecosystems today use automated package managers (e.g., APT, NPM, Maven, Cargo, and Spack), which analyze compatibility constraints among packages and select a consistent set of versions to install. Simply selecting a *compatible* version configuration is known to be NP-complete [17, 37], but there may be many valid configurations. Selecting the *best* of these requires that we solve an NP-hard constraint problem *along with* an NP-hard optimization problem [6, 41, 52].

There are many reasons to choose software versions carefully. Recent exploits in package ecosystems have highlighted issues with the strategies employed by most package managers [9] [14]. In particular, many systems will choose the latest possible version, assuming that it has the latest fixes and security updates. However, using the latest version can also break package builds. We would like to be able to balance these concerns, but threading the needle between the need to track updates and choosing known stable builds is extremely hard. The alternatives are to use *wisdom of the crowds* (choice of the majority) [42] when selecting package versions or to select the *lowest* allowed version. However, none of these prevent or detect conflict defects [7].

In the scientific software and ML ecosystems, the pain of build errors is particularly acute. Code is written in compiled languages such as C, C++, and Fortran, along with interpreted front-ends in Python and Lua. There are many compilers, package versions, language versions, and interoperability concerns. In these communities, it is common to port software to new architectures and platforms with new compatibility concerns. Existing dependency version constraints won't help here, as they are often dealing with bleeding-edge versions of packages, where compatibility information is not yet known. Unlike using a standard Linux package manager, where we are choosing among pre-existing builds that are known to work, we are choosing new configurations that have not been tested. We would like our packaging tools to help by visiting the bleeding edge to find bad configurations before we encounter them. Ideally, we could learn from related builds and increase the likelihood that a build we have never tried before will work. We call such builds "high-fidelity builds".

In this paper, we present an active-learning-based method to select configurations with high-fidelity, and thus a high likelihood of building successfully. Our approach uses adaptive sampling, a technique where samples are collected in an iterative fashion, to reduce the number of samples used to identify high-fidelity built configurations. We leverage the inherent relationship between sub-packages in the form of a dependency graph and use it to design a surrogate model to predict whether a configuration will successfully compile. Our surrogate model is a probabilistic model that gives a

\*These authors contributed equally to this work.

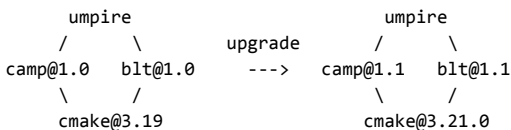
joint distribution over the dependencies to calculate a score that indicates whether a configuration will build. This provides a set of promising samples which are used by the adaptive sampling algorithm to evaluate their true objective function by running the build process. This approach will enable users to select configurations that are highly likely to build using limited evaluations, reducing the user effort and resource overhead. Our main contributions are:

- *Reliabuild*, a novel active-learning-based approach that selects high-fidelity build configurations.
- Two probabilistic surrogate models that assign build scores using: 1) wisdom of the crowd, 2) a Bayesian model that learns from dependencies among the packages.
- Automatically deriving dependency version constraints by analyzing the relative importance of different packages in the learned model.
- A detailed experimental evaluation of our approach over 31,186 builds from the E4S scientific software ecosystem.

In our evaluation, we find that *Reliabuild* selects high-fidelity build configurations automatically using far fewer samples than a randomized selection. The Bayesian as well as the wisdom of the crowd surrogate models are able to select good build configurations. However, the model based on Bayesian optimization provides better prediction accuracy. Further, *Reliabuild* provides several insights that can be used by package managers as well as developers to determine possible version conflicts as well as understand the sensitivity of the build outcome to certain packages versions.

## 2 OVERVIEW

Here is an example showcasing how a change in the requirement of one package leads to cascading issues. We look at *umpire*, a memory allocation package, which has 15 dependencies, 3 of which are shown in Figure 1.



**Figure 1: An upgrade scenario of the *umpire* memory allocation package.**

Suppose further that:

- *blt* 1.0 requires *cmake*  $\geq 3.18$ , but is incompatible with *cmake* 3.21.0 due to an as-yet unknown bug
- *camp* 1.0 depends on *cmake* 3.19 or higher, but *camp* 1.1 depends on *cmake* 3.21 or higher
- The *umpire* developers want to use *camp* 1.1 for its new features
- Upgrading *camp* pushes *cmake* to the latest 3.21.0, but this is incompatible with *blt*, and the incompatibility will cause the build to fail.
- We need to use *blt@1.1* to make this work.

To find the incompatibility, developers of *umpire* would have to try building versions of *cmake* until they find one that works

with both *blt* and *camp*. While doing so they may discover other incompatibilities. This becomes burdensome when dealing with hundreds of dependencies requiring careful evaluation of compatibility. Moreover, this would have to be done manually on different architecture and compilers to ensure they work as expected.

*Reliabuild* automates this search and makes recommendations based on what it learns. In this example, the developers of *umpire* can use *Reliabuild* adaptive sampling (1) to visit the bleeding edge to evaluate sample configurations provided by the model (Eq 3) and use those samples to update the model (Eq 1 and 4). This increases the likelihood of selecting a good build that we have never tried before.

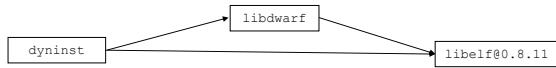
## 3 BACKGROUND

The appeal of component software is clear; dividing software into logical components enables developers to separate concerns and reuse software more easily. Experts on one part of a system can independently develop a component package and expose an interface for other teams to rely on. Teams that rely on an external dependency can easily upgrade to newer versions with bug fixes and security updates, provided that the dependency’s interface is still compatible. Ensuring this compatibility, or at least reliably predicting it, is the challenge we aim to address in this paper.

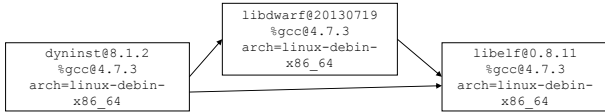
### 3.1 Versioning

Package managers rely on version metadata to determine package compatibility. Ultimately, this metadata comes from humans (package developers and maintainers) with some knowledge of the package and its dependency relationships. Developers can declare dependencies using a *fixed* version (e.g. use only version 1.9.2) or a *version range* (e.g use a version which is at least as recent as 1.0.2 or below version 1.9.2 or between two 1.8.2 – 1.9.2). While *fixed* versions improve build determinism, they limit the flexibility of a package. Supporting only a specific version can conflict with the versions required by other dependent packages, and this limits composability by making it very difficult to integrate packages into large projects. The problem with flexible version ranges is that maintainers cannot test all possible configurations. Rather than exhaustively testing the versions in a range, developers more often simply assume that compatibility is consistent across all of them.

**3.1.1 Balancing versions.** Maintainers must strike a delicate balance. They need to keep their dependency specifications as flexible as possible to benefit from bug fixes and new features of dependent packages. However, they must also avoid including incompatible updates. At the same time, package developers must be careful not to release new versions with breaking changes. Conventions such as semantic versioning, or *semver* [49], can help to identify new releases that either preserve or break compatibility, but semantic versioning relies heavily on developers to know the rules of compatibility and versioning their package releases. Such rules are complex and typically not fully understood [18] and, although developers are getting better with them, there is still a lot missing [15, 19]. Tools that identify incompatible changes through static analysis are still incomplete [32]. While testing [12] provides some coverage, it cannot guarantee [29] compatibility. Even minor version changes



(a) Abstract Spec: DAG representation of the spec: “*dyninst libelf@0.8.11*”. The spec allows freedom to the Spack concretizer to select any version for *dyninst* and *libdwarf* but constrains the *libelf* package to use the fixed version *0.8.11*



(b) Spack provides as input to the concretizer the abstract spec depicted in Figure 2a. The concretizer produces as output a concrete spec. Namely, in a concrete spec each node in the dag has all the information fully defined.

Figure 2: Abstract and concrete *spec* DAG representations.

frequently break builds, and finding a working configuration can involve manual experimentation.

**3.1.2 Version exploits.** Version flexibility has been exploited repeatedly. Most recently in the NPM ecosystem, a maintainer intentionally sabotaged a popular NPM package by removing it from NPM and deleting its repositories [14]. After it was reinstated by the community and NPM, the maintainer uploaded a version *with a higher version number* introducing changes designed to break any software that used the package. Dependent packages immediately started to fail, resulting in thousands of bug reports. Another exploit called *dependency confusion* [9] takes advantage of version priority. A project with private dependencies meant to be fetched internally may also rely on a public package repository for open source dependencies. Attackers can add an identically named, higher-versioned malicious package to an *external* repository and override the internal package, allowing them to run code on private systems.

## 3.2 Spack

In this paper we use the Spack [23] package manager to explore the combinatorial build space of packages and to understand what factors contribute to a successful build. We chose Spack for several reasons. First, because it comes from the world of High Performance Computing (HPC), Spack is designed to support building packages from source. Second, as HPC engineers frequently need to port packages to new systems, Spack is *very* flexible about versions, build parameters, and dependency configurations. It exposes parameters to users as adjustable knobs and allows a single package to be built in many different ways. This is ideal for experimentation. Finally, like system-level Linux package managers, Spack is language-agnostic. Its focus is on HPC, scientific computing, and machine learning, so C, C++, Fortran, Python, and various parallel programming models comprise most of the > 6,000 packages in its ecosystem. Scientific software is notoriously difficult to build [20, 23, 30, 35], and our aim is to simplify build configuration.

**3.2.1 Package DSL:** Spack packages are written in a domain-specific language embedded in Python, as shown in the example

```
class Kripke(CMakePackage):
    """Kripke is a 3D Sn deterministic particle transport mini-app."""
    homepage = "https://computation.llnl.gov/kripke"
    url = "https://computation.llnl.gov/kripke-openmp-1.1.tar.gz"

    version('1.2.3', sha256='3f7f2eef0d1ba5825780d626741eb0b3f026...')
    version('1.2.2', sha256='eaf9ddf562416974157b34d00c3a1c880fc5...')
    version('1.1', sha256='232d74072fc7b848fa2adc8a1bc839ae8fb5...')

    variant('mpi', default=True, description='Build with MPI.')
    variant('openmp', default=True, description='Enable OpenMP.')

    depends_on('raja', when='@1.2:')
    depends_on('mpi', when='+mpi')
    depends_on('cmake@3.0:', type='build')

    def cmake_args(self):
        return [
            '-DRAJA_PREFIX=%s' % (self.spec['raja'].prefix),
            '-DENABLE_OPENMP=%s' % ('+openmp' in self.spec),
            '-DENABLE_MPI=%s' % ('+mpi' in self.spec),
        ]
```

Figure 3: Example of a parameterized Spack package.

Kripke package in Figure 3. Each package is a class containing *directives* at the class-level to constrain the build space and *functions* that describe how to build. In the figure, version directives describe available version tarballs and their checksums, variant directives expose optional build parameters, and depends\_on directives define dependencies on other packages. Packages may require particular versions of dependencies, e.g., kripke requires cmake at version 3.0 or higher. Dependencies can also be conditional on versions, variants, or other properties. Here, kripke only depends on raja when it is at version 1.2 or higher, and it only depends on mpi when the mpi variant is enabled. The raja, mpi, and cmake packages are described by separate package templates like this one.

Directives are supplied by package maintainers to define the combinatorial space of configuration that are *possible*. Spack selects a concrete configuration, or *spec*, from this space and builds it using the recipe described in the package’s functions. Here, the package extends a base CMakePackage class that has the bulk of the standard recipe for handling the CMake build system, and the cmake\_args method tells the package how to translate the spec configuration to arguments for the cmake command that configures the build.

**3.2.2 spec syntax:** The *spec syntax* in Spack allows users to specify their own constraints on top of those provided by maintainers in package recipes. For example, often a user may only care to build a particular package with no additional constraints, or build a package with a specific compiler. Table 1 presents some examples of the spack install command, the spec syntax and a description. The spec syntax is fully recursive, in that any dependency can be constrained, just as the root node can.

**3.2.3 Concretization:** The specs in the table 1 are *abstract* – that is, they do not specify *all* aspects of build configuration. Typically, users care about a very small set of constraints, and rely on Spack to decide the rest of the build configuration. In Spack, the process of selecting a consistent build configuration is called *concretization*. In other package managers, the analogous process is typically called *dependency resolution*. Figure 2a presents an abstract spec DAG along with its corresponding *concrete* spec. The concrete graph is created by combining constraints from packages, command line, and user preferences and solving for what nodes and configuration

**Table 1: Examples of the spack install command using the spec to constrain the combinatorial build space.**

Examples	Meaning
spack install kripke	Install kripke but do not constrain the installation; Spack decides how to build.
spack install kripke %gcc@8.3.1	Build kripke with gcc 8.3.1
spack install kripke@1.2.2~mpi ^raja@1.9	Build kripke version 1.2.2 with mpi disabled and with raja version 1.9.

options should be present. Constraints not defined in the abstract spec represent degrees of freedom.

Spack’s *concretizer* is a combinatorial logic solver implemented using Answer Set Programming (ASP) [25]. It translates inputs from i) the *package-DSL*, ii) the *spec-syntax*, and iii) configuration file into a Prolog-like program that finds a dependency graph satisfying all constraints from the *package DSL* and the *spec syntax* on the command line. The concretizer eases the task of exploration because it ensures that only valid configurations are resolved. However, there can exist multiple configurations which satisfy those constraints. Typically, users must explore this space manually by specifying constraints, but our goal in this work is to automatically find good configurations *within* the set of valid configurations.

## 4 HIGH-FIDELITY BUILD CONFIGURATION SELECTION

Identifying compatible versions of a package’s dependencies is a time consuming and error-prone endeavour. Currently, the process of identifying a package configuration that successfully builds relies on individuals with in-depth knowledge of various packages and their interactions with their dependencies. In the absence of such knowledge users have to resort to exploring the build configuration space by trying out various combinations of packages and their dependencies. The space formed by all possible options for a package and dependencies is immense, and as a result an exhaustive search of that space is impractical. One can resort to random sampling, but that is still a tedious process, and if system parameters change it would have to be redone all over again. Ideally, we want an automated process that can select configurations that are highly likely to build based on a smaller set of samples.

The trick to learning quickly is to pick samples carefully, and we develop an active-learning-based approach to identifying high-fidelity package build configurations using only a limited set of samples. Active learning algorithms select data from which to learn in order to achieve a specific objective. This is especially suitable when the true objective function evaluations are expensive, as is in the case of building packages. Below, we describe the problem formally and provide details about the iterative sampling, design of the model, and describe the selection algorithm in its entirety.

Let the root package we are building require  $n$  dependencies, each of which is represented by  $X_i$ ,  $i \in 1, \dots, n$ . A package  $X_i$  can take on values  $x_i \in v_{i,1}, \dots, v_{i,m_i}$ . We use  $\mathbf{x} \equiv [x_1, \dots, x_n] \in \mathcal{X}$  to represent a configuration of the  $n$  dependent packages as a vector.

Let  $f : \mathcal{X} \rightarrow \{0, 1\}$  be a function representing the outcome of building a configuration. If a configuration  $\mathbf{x}$  builds successfully, then  $f(\mathbf{x}) = 1$  and if it fails to build, then  $f(\mathbf{x}) = 0$ . Our goal is to find *some* configuration  $\mathbf{x}^* \in \mathcal{X}$ , from the space of all possible configurations  $\mathcal{X}$ , that successfully builds. This can be represented by the following objective function.

$$\mathbf{x}^* = \arg \max_{\mathbf{x}} f(\mathbf{x})$$

Note that, while there may be many configurations that successfully build, we are interested in finding any one of those.

### 4.1 Iterative Configuration Selection Algorithm

We use an iterative process for sampling configurations that are likely to build. We construct a probabilistic surrogate model that predicts the probability of a configuration to compile successfully. This model is used to sample a set of highly likely candidates which are then built. The outcome is then used to update the model and the process is repeated. Thus, our method follows an adaptive iterative scheme, where it alternates between updating the model and using it to make choices about which configurations to investigate next. Our method is sketched in Algorithm 1.

---

#### Algorithm 1 Pseudocode for Adaptive Sampling

---

```

1: procedure SAMPLE( $f, \mathcal{M}_0, \mathcal{S}, T$ )
2:    $\mathcal{H}_0 \leftarrow \emptyset$ 
3:   for  $t = 1$  to  $T$  do
4:      $\mathbf{x}_t^* \leftarrow \arg \max_{\mathbf{x}} \mathcal{M}_{t-1}(\mathbf{x})$ 
5:      $y_t^* \leftarrow f(\mathbf{x}_t^*)$  ▷ Expensive process
6:      $\mathcal{H}_t \leftarrow \mathcal{H}_{t-1} \cup \{(\mathbf{x}_t^*, y_t^*)\}$ 
7:      $\mathcal{M}_t \leftarrow \text{FITMODEL}(\mathcal{M}_{t-1}, \mathcal{H}_t)$ 
8:   end for
9:   return  $\mathcal{H}_T$ 
10: end procedure

```

---

In each iteration  $t$ , a single configuration  $\mathbf{x}_t^*$  is selected, from a list of candidate configurations. As evaluating the true objective function for all the candidate configurations is expensive, we use a surrogate model  $\mathcal{M}_{t-1}$  to score the candidates. The surrogate model is constructed from the observation history,  $\mathcal{H}_{t-1}$ , and is cheap to evaluate. The candidate  $\mathbf{x}_t^*$  with the maximum score value is selected at each iteration and is used to evaluate the true objective function,  $y_t^* = f(\mathbf{x}_t^*)$ , by running the build process for the configuration  $\mathbf{x}_t^*$ . The result of the evaluation is then added to the observation history  $\mathcal{H}_{t-1}$  to yield the updated history  $\mathcal{H}_t$ . Finally,  $\mathcal{H}_t$  is used to update the model, yielding  $\mathcal{M}_t$ , and choose the next configuration  $\mathbf{x}_{t+1}^*$  at iteration  $t + 1$ . This process continues iteratively until iteration  $t = T$ . The algorithm starts off with a small set of 20 initial samples drawn uniformly at random. All of them are evaluated to obtain initial observation history  $H_0$  and then create the initial surrogate model  $\mathcal{M}_0$ . Based on our evaluation, the model performance is insensitive to the initial sample size. We observed a drop in the performance on both extremes (too many or too few initial samples). For example, when bootstrapped with a larger fraction of samples (>=50%), only a smaller fraction of the samples are

selected by the model and the rest are selected uniformly at random. In contrast, bootstrapping with too few samples, such as 1, does not capture correctly the distribution of the package versions therefore resulting in poor modeling behavior. This method is a form of Sequential Model Based Global-Optimization methods (SMBO), and has been used for black box optimization of functions that are expensive to evaluate [31, 33, 54].

## 4.2 Probabilistic Surrogate Model

A surrogate model is a model that approximates the true objective function,  $y = f(\mathbf{x})$ , and provides a significantly cheaper method for computing the approximate objective for all samples. Given a space  $\mathcal{X}$  of all possible build configurations, our goal is to design a surrogate model that provides a score indicating whether a configuration  $\mathbf{x} \in \mathcal{X}$  is likely to build. A good surrogate model will assign a high score to each good configuration. As the surrogate is significantly cheaper to evaluate than the true objective function, the iterative sampling algorithm (see section 4.1) can use it to rank all the configurations and select the most promising configuration whose true objective function can then be evaluated by running the build process. We describe two surrogate models that are incorporated in *Reliabuild*.

**4.2.1 Bayesian Model.** To construct this model, we use the algorithm used in Bayesian optimization [43], which approximates the objective function with a surrogate objective that is cheaper to compute. Bayesian optimization has been used for tuning the hyperparameters of deep neural networks [8], auto-tuning parameters of HPC systems and applications [5, 11, 40], where training is a very expensive process. The surrogate model  $\mathcal{M}(\mathbf{x})$  computes Expected Improvement ( $\mathcal{I}$ ) [33] using a probabilistic model  $p_{y|\mathbf{x}}(y | \mathbf{x})$  for some configuration  $\mathbf{x}$ . The Expected Improvement  $\mathcal{I}$  (eq. (3)) is the expected margin by which the true objective  $f(\mathbf{x})$  will be 1 (successfully builds).

In the Bayesian optimization method, Bayes rule is used to define  $p_{y|\mathbf{x}}(y | \mathbf{x})$  in terms of  $p_{\mathbf{x}|y}(\mathbf{x} | y)$ , yielding:

$$p_{y|\mathbf{x}}(y | \mathbf{x}) = \frac{p_{\mathbf{x}|y}(\mathbf{x} | y)p_y(y)}{p_{\mathbf{x}}(\mathbf{x})} \quad (1)$$

Here,  $p_y(y)$  is the prior distribution for  $y$  and  $p_{\mathbf{x}}(\mathbf{x}) = \int p_{\mathbf{x}|y}(\mathbf{x} | y)p_y(y)dy$  is the marginal distribution of  $\mathbf{x}$ .

To model  $p_{\mathbf{x}|y}(\mathbf{x} | y)$  we split  $y$  into two possibilities: 1) successfully builds i.e.  $y = 1$ , and 2) fails to build i.e.  $y = 0$ . This enables us to define the probability distribution  $p_{\mathbf{x}|y}(\mathbf{x} | y)$  in terms of two probability density functions;  $p_g(\mathbf{x})$  for good configurations, and  $p_b(\mathbf{x})$  for bad configurations:

$$p_{\mathbf{x}|y}(\mathbf{x} | y) = \begin{cases} p_g(\mathbf{x}) & \text{if } y = 1 \\ p_b(\mathbf{x}) & \text{if } y = 0, \end{cases} \quad (2)$$

Finally, the Expected Improvement  $\mathcal{I}(\mathbf{x})$  can now be written as:

$$\mathcal{I}(\mathbf{x}) = \frac{1}{\alpha + \frac{p_b(\mathbf{x})}{p_g(\mathbf{x})}(1 - \alpha)} \quad (3)$$

We can compute  $\mathcal{I}(\mathbf{x})$  for each candidate configuration and choose the one with the highest value as the candidate  $\mathbf{x}_t^*$  that is most likely to build.  $\mathbf{x}_t^*$  can then be used for performing the full evaluation  $y_t^* = f(\mathbf{x}_t^*)$ .

To construct  $p_g(\mathbf{x})$  and  $p_b(\mathbf{x})$ , we leverage the underlying dependency structure of the package. We represent a package and its dependencies using an undirected graph, and consider a joint distribution over the dependencies to calculate the score. Let  $G = (X, E)$  denote the undirected graph corresponding to the root package of interest, where  $X$  is the set containing the package and the dependencies and  $E$  is the set of edges as given by the dependency structure of the package.

Our surrogate model uses a probability density defined over the prediction values  $y$  given  $\mathbf{x}$ . Estimating the full joint distribution over the parameter space is not feasible. However, we can leverage the inherent relationship between the packages provided by the package dependency graph to obtain a factorized distribution. This allows us to represent the joint probability distributions as a product of factors with smaller number of variables. Let the factors be  $\{F_j(V_j)\}_m$ , where  $V_j \subseteq \{X_1, \dots, X_n\}$  is a subset of the variables. Then we can write the joint distribution  $p_g(\mathbf{x})$  as

$$p_g(\mathbf{x}) = \prod_{j=1}^m F_j^{(g)}(V_j) \quad (4)$$

We use the  $(g)$  superscript to denote the factors corresponding to distribution for the good configurations  $p_g$ . We consider factors of at most size two, that are defined for every node and edge of the dependency graph. For example, a factor  $F^{(g)}(\{X_j, X_k\})$  corresponding to an edge between the node  $j$  and the node  $k$  captures the likelihood of that pair compiling.

The underlying assumption behind this factorization is that the probability of a package building will depend on whether that package and its dependencies can be built successfully. In the event where a package and its immediate dependency fails to build, the package fails to build because its requirements cannot be satisfied. Hence, instead of estimating the full joint distribution, we can reasonably approximate the probability of a package building as a factorization over the probability distribution of the pairwise distributions between each parent and child (given by the dependency graph).  $p_b(\mathbf{x})$  is constructed similarly to Eq 4 using the same factorization  $\{V_j\}_m$  but a different set of factor functions  $F_j^{(b)}$ .

**4.2.2 Wisdom of the Crowd.** This approach uses the majority opinion to select versions for packages [42]. The probability of selecting a version for a package is directly proportional to the number of times it occurs in successfully built configurations in the observed data. This can be represented as a probability distribution as shown below.

$$p(\mathbf{x}) = \prod_{i=1}^n p_g(v_i) \quad (5)$$

where  $p_g(v_i)$  is the probability of finding package  $v_i$  of  $X_i$  in the observed set of good configurations.

## 4.3 Pairwise Importance Analysis

Packages can influence the build outcome differently. Some packages and package pairs might be compatible across all possible versions and some might have very strict compatibility rules. Identifying the packages or pairs of packages that are crucial for the success or the failure of the build process enables users to better

understand the most likely sources of failure. We use the two probability densities from the surrogate model to derive a metric that tells us the relative significance of each package in build outcomes. Recall that the surrogate model maintains two distributions: a distribution of package versions based on high-fidelity configurations  $p_g$ , and a distribution of package versions based on bad configurations  $p_b$ . The degree to which these distributions differ from one another is an indication of how sensitive the build outcome is to the version of this package. We use Jensen-Shannon (JS) divergence to compute the difference between the two distributions.

For two probability distributions  $P$  and  $Q$  defined in the same probability space  $X$ , the JS divergence is defined using  $M = (P + Q)/2$  as

$$D_{JS}(P, Q) = \frac{1}{2}D_{KL}(P, M) + \frac{1}{2}D_{KL}(Q, M) \quad (6)$$

$$D_{KL}(P, M) = \sum_{x \in X} P(x) \log \frac{P(x)}{M(x)} \quad (7)$$

Here  $D_{KL}(P, M)$  is the Kullback–Leibler (KL) divergence from  $M$  to  $P$ . We use this metric to study the influence of different packages and pairs of packages on the build outcomes.

## 5 EXPERIMENTAL SETUP

In this section we provide details about the dataset used, data collection mechanism, and the metrics used for our evaluation.

### 5.1 Evaluation Dataset

We evaluated our model on several packages from the Extreme-scale Scientific Software Stack (E4S). E4S provides open source software packages for developing, deploying and running scientific applications on high-performance computing (HPC) platforms. These software packages are implemented in different programming languages such as C/C++, FORTRAN, Python, Lua, and others. E4S uses Spack for managing software packages. Table 2 shows the list of packages from E4S that were used for our evaluation.

### 5.2 Data Collection

We explored the build space for several Spack packages from the E4S software stack. The combinatorial package space consists of millions of configurations. So, to create the dataset we randomly select a set of concrete specs from that space. We refer to these concrete specs as configurations. The selection samples only different versions and does not sample variants or compiler options.

We evaluate whether each configuration successfully builds or not. Initially, we create a DAG in which the nodes are the unique packages of all the sampled configurations. Common dependencies among configurations appear only once in the DAG and, thus, we built those only once. We implement a parallel build process of the entire DAG in a distributed system. The process follows a *farmer-worker* parallel paradigm. A single node operates as the farmer. The farmer traverses the DAG and selects the nodes ready for installation. A node is ready when all the dependencies of the node have been successfully built. Once the farmer gathers these nodes, it assigns their installation to workers. The workers operate in parallel within the distributed cluster and issue the

spack install "specification" command to build the individual nodes. Once all workers finish their installations, the farmer identifies which installations failed or succeeded and updates the DAG. When a dependency fails, all dependent nodes are marked recursively as failed. The process continues until all nodes in the DAG are marked as failed or successful. In the end, the farmer exports the DAG information and the node installation status to a pandas [39] *DataFrame* and stores it into permanent storage.

The farmer-worker installation protocol relies on Spack being parallelism aware and allowing multiple Spack instances to execute concurrently without corrupting internal Spack data structures. Spack implements a locking mechanism through the network filesystem protocols and uses the lock when updating internal data structures (such as the user spack database). During our evaluation, we repeatedly pushed the locking mechanism to its limits. The process exposed concurrency bugs in the mutual exclusion algorithm implemented in Spack. We communicated and fixed those bugs resulting in a significantly improved locking mechanism. In the end, the farmer-worker protocol coupled with the improved locking mechanism resulted in executing up to 432 concurrent spack instances with each instance using 9 CPU cores.

In Table 2 we show the total number of configurations for each package, the number of successful builds, and the number of dependency packages. Naively, building this data set would require 83,796,782 package installations, but Spack identifies unique builds with a Merkle hash of their configuration metadata. That is, the uniqueness of a *configuration* depends on the root's metadata and the metadata of all of its dependencies. By hash, there were 56,645 unique package configurations (including root packages and dependencies). We skipped the installation of 19,662 nodes due to a failed dependency build, and we performed 31,186 total package installations out of which 28,157 were successful. This dataset was collected on a 3,018-node Intel Xeon cluster with 36 cores per node.

### 5.3 Metrics for Evaluation

A good sampling selection algorithm is expected to identify high fidelity configurations while observing fewer samples. As a result we use the following metrics to evaluate our approach.

**Precision ( $\mathcal{P}$ )** tells us what fraction of the samples contain configurations that built successfully.

$$\mathcal{P}(\mathcal{H}) = \frac{|\{x|x \in \mathcal{H}, f(x) = 1\}|}{|\mathcal{H}|} \quad (8)$$

**Recall ( $\mathcal{R}$ )** gives the ratio of the configurations that successfully built and included in the sampled set to the actual number of configurations that successfully built in the entire sample space.

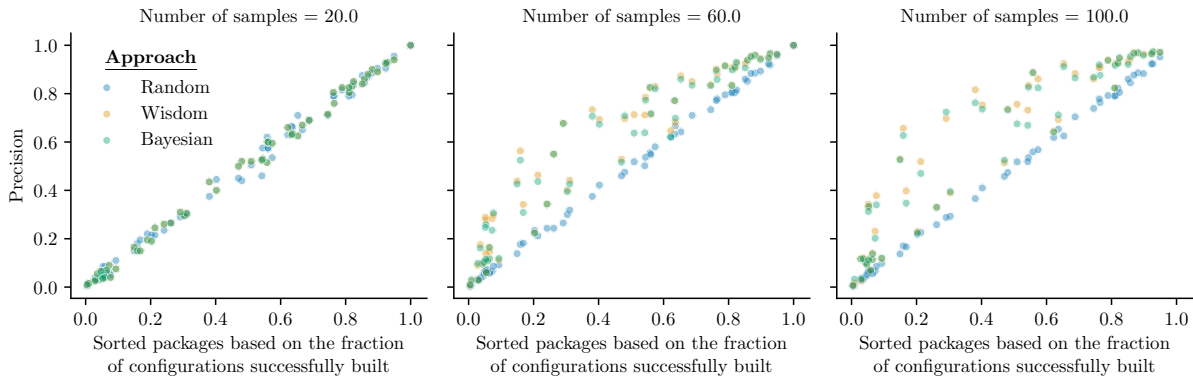
$$\mathcal{R}(\mathcal{H}) = \frac{|\{x|x \in \mathcal{H}, f(x) = 1\}|}{|\{x|\forall x, f(x) = 1\}|} \quad (9)$$

**Area Under the Precision-Recall Curve (AUPRC):** A PR curve shows the trade-off between precision ( $\mathcal{P}$ ) and recall ( $\mathcal{R}$ ) across different sample sizes. AUPRC is obtained by calculating the area under the PR curve, and is calculated as shown below.

$$AUPRC = \sum_{k=1}^N \mathcal{P}(\mathcal{H}_k) \Delta \mathcal{R}(\mathcal{H}_k) \quad (10)$$

**Table 2: The number of tested configuration for various packages within the E4S software ecosystem. In total we explore 31, 186 unique root package builds resulting in examining 83, 796, 782 root packages and dependencies.**

Package	Configs(#)	Good (#)	Deps (#)	Package	Configs(#)	Good (#)	Deps (#)	Package	Configs(#)	Good (#)	Deps (#)
abyss	892	133	36	adios	58	11	48	amrex	616	543	40
ascent	320	14	32	axom	166	5	40	bolt	996	53	22
cabana	428	130	42	caliper	289	221	41	chai	298	19	16
conduit	599	241	29	darshan-runtime	190	163	39	hypre	322	264	40
hpx	799	384	45	heftte	292	163	33	hdf5	691	641	40
gmp	636	18	19	globalarrays	323	262	40	fortrilinos	146	8	30
faodel	361	269	29	kokkos	625	294	13	kokkos-kernels	567	115	14
libnrm	923	48	40	mercury	990	539	49	metall	275	261	14
mfem	441	348	119	mpark-variant	198	126	13	ninja	621	503	23
omega-h	999	213	42	openmpi	333	12	114	openpmd-api	90	28	41
papyrus	977	672	40	parallel-netcdf	345	25	39	petsc	403	26	141
phist	852	434	183	plasma	993	378	14	pumi	990	914	40
py-libensemble	381	18	56	py-petsc4py	119	7	50	qt	670	6	177
qthreads	94	94	22	qwt	87	21	102	raja	82	46	15
rempi	848	135	39	scr	30	20	65	slate	169	97	44
slepc	126	33	45	stc	692	3	41	sundials	708	636	40
superlu-dist	624	58	43	swig	215	183	18	sz	652	497	14
tasmanian	860	709	40	trilinos	1000	168	65	turbine	565	369	34
umap	702	610	13	umpire	82	52	15	unifyfs	246	13	41
upcxx	994	539	34	variorum	989	288	24	veloc	667	51	37
zfp	540	336	13								



**Figure 4: Precision for all the packages. Each point on these scatter plots corresponds to a package. Both the models of *Reliabilbuild* (Bayesian and Wisdom of the crowd) have a significantly higher precision than *Random Selection*. *Reliabilbuild* is bootstrapped with 20 samples drawn uniformly at random, so the performance of all the approaches is similar for sample size 20.**

where  $N$  is the total number of configurations in the collection,  $\mathcal{P}(\mathcal{H}_k)$  is the precision for  $k$  samples, and  $\Delta\mathcal{R}(\mathcal{H}_k)$  is the change in recall that happened between  $k - 1$  and  $k$  samples. The AUPRC is a useful metric in problem settings like this where we care about finding positive samples. In our case, AUPRC will provide a measure of how well the surrogate model is able to select high-fidelity build configurations. AUPRC value of 1 indicates a model that is able to select high-fidelity build configurations perfectly.

## 6 EVALUATION

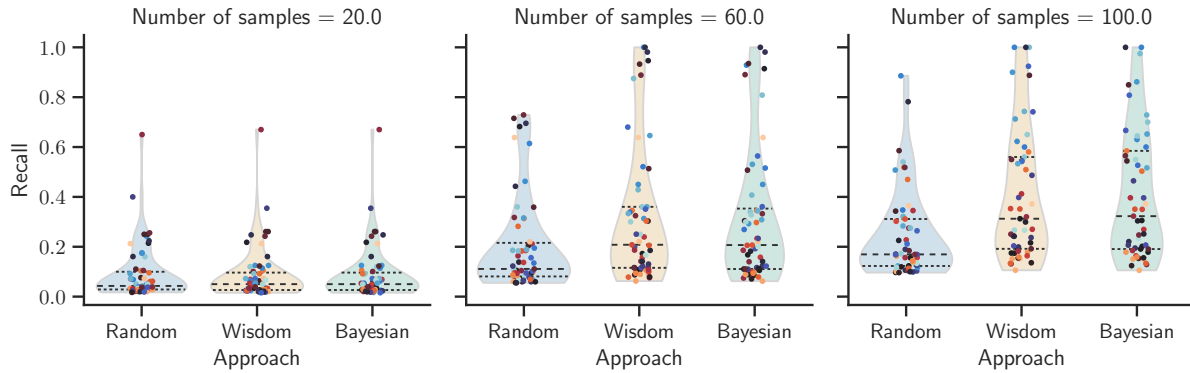
We answer three research questions by evaluating our approach on the datasets listed in Table 2. The RQ1 aims to understand whether we can use empirical information to learn a model that selects working software configurations (Section 6.1). The RQ2 aims to identify packages that are more likely to affect build compatibility

(Section 6.3). The RQ3 offers insight into different packages version compatibility (Section 6.4).

For our evaluation, we compare the performance of all the methods for a range of samples by running the algorithm 10 times and reporting the mean for each evaluation metric. We evaluate *Reliabilbuild* by comparing the Bayesian model and the Wisdom model against *Random Selection*. In *Random Selection*, the configurations are selected uniformly at random from the database.

### 6.1 Evaluation of Build Configuration Selection

A package dataset consists of several build configurations, of which only a subset successfully built. The goal of configuration selection is to identify high-fidelity builds using fewer samples. Figure 4 shows the  $\mathcal{P}$  for different packages at different sample sizes. The x-axis lists the packages sorted based on their build success rate



**Figure 5: Recall for all the packages.** The violin plots show the probability density of recall values along with the interquartile means for different sample sizes. Each point shown here corresponds to a package color-coded based on their overall build success rates (light blue to dark red corresponds to build success rates ranging from 0 to 1).  $\mathcal{R}$  value of 1 indicates that all the good configurations are included in the selected samples. Both the models of *Reliabil* (*Bayesian* and *Wisdom of the crowd*) have a higher recall score than *Random*, and for some packages attains  $\mathcal{R}$  of 1.

(number of good configurations divided by the total number of configurations). The figure shows that both *Reliabil* models (*Bayesian* and *Wisdom of the crowd*) have a significantly higher fraction of high-fidelity build configurations compared to *Random Selection*. As mentioned in section 4.1, we bootstrap our search with an initial set of 20 samples chosen uniformly at random, which is why all the algorithms have a similar performance for the sample size of 20. Figure 5 shows the  $\mathcal{R}$  for different packages at different sample sizes. The  $\mathcal{R}$  metric for *Bayesian* and *Wisdom of the crowd* is much higher than *Random Selection* indicating that these approaches identify several, if not all, high-fidelity build configurations with far fewer samples. Note that packages with a low build success rate have only a few good configurations, hence low precision values. For packages with a higher build success rate,  $\mathcal{P}$  would be close to 1, however  $\mathcal{R}$  of 1 can only be attained when at least that many samples are selected. For some packages, our approach achieves an  $\mathcal{R}$  of 1 with just 60 or 100 samples, indicating that *Reliabil* selected all the high-fidelity builds.

## 6.2 Evaluation of the Models

While the metrics  $\mathcal{P}$  and  $\mathcal{R}$  show us the ability of the models to select a sample that successfully builds, they don't tell us how well the model will perform if it has to identify all the good configurations. We use the AUPRC metric from Eq 10, which calculates the area under the Precision Recall curve, to compare the different algorithms. The ideal value for AUPRC is 1 and the higher the value the better the algorithm is at selecting good configurations. We divide the data into test and training sets, train the models on the training set and use it to evaluate the test set. Half of the dataset was used for training and it was tested on the other half. The model used the training set and iteratively selected samples and updated the model. 100 samples were selected. We used this model to evaluate the test set and calculated the AUPRC metric. Figure 6 is a violin plot showing the AUPRC metric for all packages in our dataset. It also shows the probability density of the data at different AUPRC values as well as the interquartile range. It can be seen that the *Bayesian*

model has a higher AUPRC value than the *Wisdom of the crowd* model, which indicates that *Bayesian* is better at selecting high fidelity build configurations. The figure also shows the distribution of AUPRC values for different packages, where each point represents a package color-coded based on their build success rates (light blue to dark red corresponds to build success rates ranging from 0 to 1). The packages with higher build success rates tend to have higher AUPRC values because they are relatively easier to identify and even *Random Sampling* can perform fairly well. For packages with lower build success rates (indicating fewer high-fidelity builds in the build configuration samples), *Bayesian* and *Wisdom of the crowd* is able to achieve a higher AUPRC than *Random Sampling* indicating that they are able to select high-fidelity builds in tough cases. We also show the mean AUPRC for different recall cutoffs in Figure 7. While both the models from *Reliabil* have similar performance for smaller recall cutoff, *Bayesian* performs better with a larger recall cutoff, indicating that it is able to more easily identify all the high-fidelity builds.

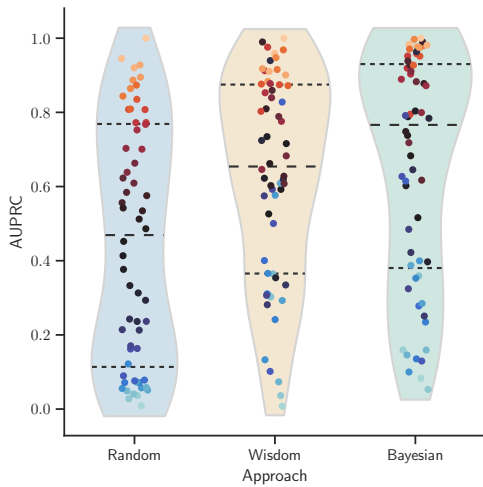
## 6.3 Package Importance Analysis

A particular choice of version for packages can significantly affect the build outcome. However, not all packages impact the application build equally. Some of them are more sensitive than others, and it would be very helpful for package managers and users to be aware of them. Table 3 shows the top five packages or parent child pairs (represented as *parent+child*, for example *autoconf+m4*) that are most sensitive and is likely to impact the outcome of the build process of the root package. The importance score calculated using Eq 7. As can be seen from Table 2, each root package contains tens if not hundreds of dependencies, and the number of package pairs formed by parent and child is on the order of hundreds. Having the knowledge about which packages influence build outcome can be used by package maintainers to exercise caution before integrating new versions of those packages. And in a rapidly evolving ecosystem guided user decisions are highly valued since it is infeasible to search the entire combinatorial space. Additionally, instead of using



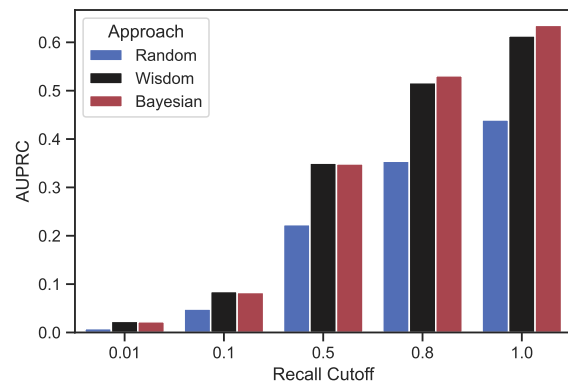
**Table 3: Relative ranking of dependencies for different root packages based on their importance. A parent child pair is represented by *parent+child* as in the case of *autoconf+m4*.**

Root package	Dependency ranking					
abyss	autoconf: 0.37	autoconf+m4: 0.37	autoconf+perl: 0.37	libtool+autoconf: 0.29	abyss+autoconf: 0.27	
adios	autoconf+m4: 0.27	autoconf+m4: 0.27	autoconf: 0.27	libtool: 0.22	libtool+m4: 0.22	
ascent	vtk-h+openmpi: 0.14	vtk-h: 0.14	vtk-h+vtk-m: 0.14	conduit+zlib: 0.12	conduit+hdf5: 0.12	
axom	lua: 0.08	lua+ncurses: 0.08	lua+readline: 0.08	lua+unzip: 0.08	axom+openmpi: 0.07	
bolt	autoconf+perl: 0.37	autoconf+m4: 0.37	autoconf: 0.37	automake+autoconf: 0.32	automake+perl: 0.30	
hypre	openblas+perl: 0.07	openblas: 0.07	hypre+openblas: 0.03	hypre+mpich: 0.02	mpich+findutils: 0.01	
hpx	hpx+boost: 0.24	hpx+hwloc: 0.24	hpx+pkgconf: 0.24	hpx+python: 0.24	hpx: 0.24	
heffte	heffte: 0.35	heffte+openmpi: 0.30	heffte+fftw: 0.24	cuda+libxml2: 0.19	mpich+findutils: 0.19	
hdf5	mpich+findutils: 0.03	mpich+pkgconf: 0.03	mpich+libxml2: 0.03	mpich: 0.03	mpich+libpciaccess: 0.03	
ninja	ninja+python: 0.03	python+ncurses: 0.01	python+readline: 0.01	python+pkgconf: 0.01	python+libffi: 0.01	
omega-h	omega-h+zlib: 0.24	trilinos: 0.24	trilinos+openblas: 0.24	omega-h: 0.24	omega-h+trilinos: 0.18	
openmpi	json-c: 0.30	mariadb+lz4: 0.30	meson: 0.30	gmp: 0.30	python+libffi: 0.30	
openmpi-api	hdf5: 0.19	hdf5+zlib: 0.19	hdf5+openmpi: 0.19	hdf5+pkgconf: 0.19	hdf5+cmake: 0.19	
papyrus	papyrus+mpich: 0.11	cmake+ncurses: 0.08	cmake: 0.08	papyrus+cmake: 0.08	mpich+findutils: 0.04	
plasma	plasma: 0.52	plasma+openblas: 0.26	openblas+perl: 0.13	openblas: 0.13	plasma+cmake: 0.12	
pumi	pumi+mpich: 0.02	mpich+findutils: 0.02	mpich+libxml2: 0.02	mpich: 0.02	mpich+libpciaccess: 0.02	
py-petsc4py	py-mpi4py+python: 0.11	hypre: 0.11	hypre+openmpi: 0.11	py-mpi4py+py-setuptools: 0.11	hypre+openblas: 0.11	
qthreads	numactl+libtool: 0.00	numactl+autoconf: 0.00	bzip2+diffutils: 0.00	util-macros: 0.00	hwloc: 0.00	
raja	raja: 0.17	raja+cmake: 0.17	blt: 0.08	blt+cmake: 0.08	raja+blt: 0.07	
rempi	autoconf+perl: 0.52	autoconf+m4: 0.52	autoconf: 0.52	rempi+autoconf: 0.50	automake+autoconf: 0.35	
scr	libyogrt: 0.08	libyogrt+slurm: 0.08	scr+libyogrt: 0.07	scr+cmake: 0.05	scr+dtcmp: 0.05	
slepc	hypre: 0.42	hypre+openmpi: 0.42	hypre+openblas: 0.42	superlu-dist+cmake: 0.39	superlu-dist: 0.39	
superlu-dist	openblas+perl: 0.19	openblas: 0.19	cmake+ncurses: 0.11	cmake: 0.11	metis+cmake: 0.11	
sz	cmake: 0.12	cmake+openssl: 0.12	cmake+ncurses: 0.12	sz+cmake: 0.04	sz+zstd: 0.03	
trilinos	trilinos: 0.19	py-setuptools+python: 0.09	py-setuptools: 0.09	trilinos+mpich: 0.08	cmake+ncurses: 0.08	
upcxx	upcxx: 0.24	upcxx+mpich: 0.22	upcxx+python: 0.11	mpich+findutils: 0.08	mpich+libxml2: 0.08	
variorum	hwloc: 0.38	hwloc+ncurses: 0.38	hwloc+libpciaccess: 0.38	hwloc+pkgconf: 0.38	hwloc+libxml2: 0.38	
veloc	veloc: 0.58	veloc+openmpi: 0.58	veloc+cmake: 0.58	veloc+openssl: 0.58	veloc+libpthread-stubs: 0.47	
zfp	zfp+cmake: 0.06	cmake: 0.04	cmake+openssl: 0.04	cmake+ncurses: 0.04	bzip2+diffutils: 0.00	



**Figure 6: AUPRC metric for all packages. Each point represents a package and the color is assigned based on their build success rate (light blue to dark red corresponds to success rates ranging from 0 to 1). This violin plot shows the probability density of the data at different AUPRC values. The *Bayesian* model has a higher AUPRC value than *Wisdom of the crowd* model indicating that *Bayesian* is a better model for identifying all the high-fidelity build configurations.**

flexible version ranges for those packages, package maintainers

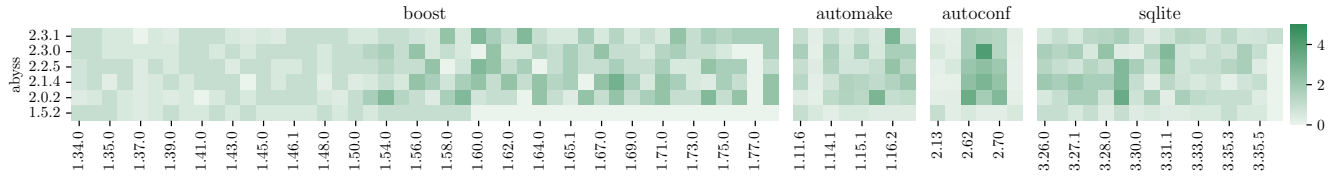


**Figure 7: Mean AUPRC metric for all packages at different recall cutoffs. While both the models from *Reliabuild* performs similar for small cutoffs, *Bayesian* performs better for larger cutoffs indicating that it is better at selecting all the high-fidelity builds.**

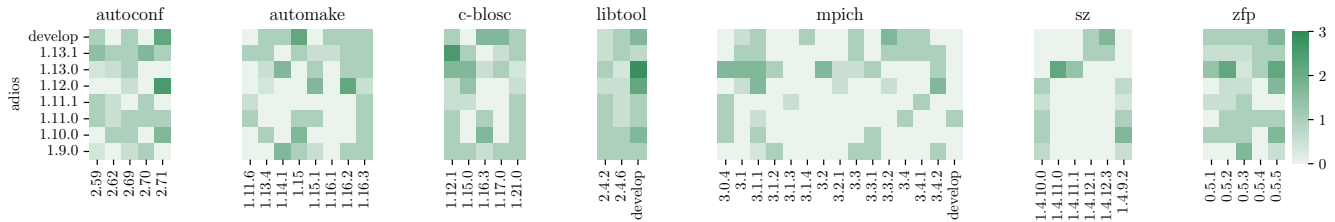
could opt for fixed versioning instead to increase the probability of getting a successful build configuration.

### 6.4 Extracting Package Dependency Constraints

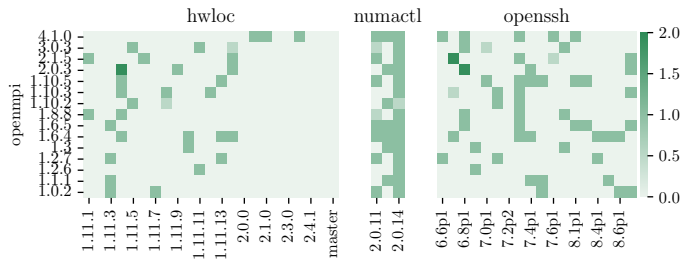
We use all the data from the dataset to build the *Bayesian* model and analyze incompatibility between various packages. We utilize



**Figure 8: Heatmap of Abyss and its dependencies with scores indicating which version pairs are highly likely to build. It shows that a newer version of Abyss (version greater than 1.5.2) is not likely to build with an older version of Boost (version 1.60.0).**



**Figure 9: Heatmap of Adios and its dependencies with scores indicating which version pairs are highly likely to build.**



**Figure 10: Heatmap of OMPI and its dependencies with scores indicating which version pairs are highly likely to build.**

the good and the bad probability densities from Eq 2 pertaining to a parent child pair and calculate the  $\mathcal{EI}$  from Eq 3 for the different package version combinations. Figures 8, 9, 10 show the heatmap indicating which version pairs are highly likely to build and which ones are not. A lower score indicates that the pair is not likely to build. This analysis was done for all the packages, but in the interest of space we show only a few of them here. Figure 8 shows which version pairs are incompatible in the case of Abyss. One of the insights provided by the figure is that a newer version of Abyss (version greater than 1.5.2) is not likely to build with an older version of Boost (version 1.60.0). Similarly, we can see from the figure that an older version of autoconf, such as 2.13, is most likely incompatible with versions of Abyss greater than 2.0.2. Moreover, note that the version pair combinations that are highly likely to build is a much smaller set which indicates that this package pair is highly sensitive to version changes. As a result, autoconf is among the top 5 sensitive packages for Abyss as shown in Table 3. Similar analysis can be done for Adios and OpenMPI based on the Figures 9 and 10. This information can be leveraged by the package manager to introduce new package dependency constraints so as to avoid configurations that can result in build failures.

## 7 THREATS TO VALIDITY

**External validity:** Reliabuild can be applied to package managers that support compiling on the target system, but it isn't really applicable when the package manager provides already built packages such as Maven and npm. We are specifically studying builds. Our study closely resembles and can generalize to ecosystems such as Debian, RHEL, and other Linux distros.

**Transferability:** While the adaptive sampling algorithm should transfer, due to differences in constraints the outcome of the compatibility analysis might not. However, we can transfer the learning from another architecture by bootstrapping with the model from a source architecture and sampling fewer samples from the target architecture.

**Credibility:** To validate our active-learning-based framework, we evaluated it on real-world problems by running on 61 packages from the Extreme-scale Scientific Software Stack (E4S) ecosystem. The precision, recall values were calculated by running it ten times and taking the mean.

## 8 RELATED WORK

In 1997 the notion of “software release management” [53] was introduced for large collections of independent packages. During that period Linux distributions broadly adopted the [21, 27] package managers. The version selection problem is NP-complete and can be encoded as SAT and Constraint Programming (CP) problems [17, 37]. Since then the idea of customizable solvers provide modular package managers [4]. The work focuses on the Common Upgradeability Description Format (CUDF). Commonly, this file format is used by the front-end of package managers to describe upgradeability scenarios to the back end-solvers. Multiple implementations are proposed to solve the CUDF scenario which employ Mixed-Integer Linear Programming, Boolean Optimization, and Answer Set Programming [6, 24, 41]. Such solutions have been adopted by various Linux distributions [3]. Today, complete solvers

are being broadly adopted by various package managers. For example, PIP recently switched to a new proper solver [50]. Dart now uses CDCL SAT solver called PubGrub [55], and Rust's Cargo [1] package manager is moving towards this approach [2].

Despite the completeness of the modern package manager solvers, the package managers rely heavily on a correct pre-selection of version dependencies among packages and detection of possible conflicts. A popular solution is semantic versioning (*semver*) [49]. In semantic version *semver* a version number characterize package compatibility and breaking changes. Semantic versioning heavily relies on developers knowing the rules of compatibility and versioning their package releases correctly. Reports and empirical studies [15, 19] indicate a broad adoption of *semver*. However, the rules are complex, and thus not fully understood [18], and therefore frequently resulting into *semver* misuse [47]. In the end, developers drop the usage of semantic versioning and use fixed versioning disregarding all the advantages of flexible versioning schemes. More importantly though, as stated in [22] if all projects, in the end, follow a fixed versioning approach conflicts will eventually arise that prevent all packages from building.

In [34] the authors study how frequently and how soon package maintainers update their dependencies to include the latest release of dependent libraries. Interestingly, package maintainers rarely update their dependencies, resulting in software systems which contain known vulnerabilities. The results indicate that a heavy reliance on libraries in contemporary projects often result in the formation of complex inter-dependency relationships inside that project. Due to such inter-dependency issues developers are reluctant to update their dependencies.

There have been multiple approaches that suggest specific dependency versions. The wisdom of the crowd [42] is a popular approach. Where maintainers depend on the most popular, or highly used, library versions. Based on such information users and maintainers can more easily select library dependencies and avoid installation errors. Other tools identify incompatibilities between versions at the binary level [13]. In contrast, some propose techniques for dependent packages to automatically perform code changes and adopt to conflicts introduced by their dependencies [56]. Multiple works [28, 44–46, 48, 51] focus on suggesting third party libraries to use in projects. In contrast to our approach, these efforts focus on suggesting new libraries, and thus new dependencies. Our work focuses on a different facet of the problem, how to select which version of an existing third-party library to use in a given configuration.

These works provide mechanisms to prevent using conflicting libraries and thus avoid performing failing installations. In our approach we embrace and accept the existence of failing build configurations. However, we provide an active learning mechanism through autotuning that learns which configurations tend to fail and after a reasonable amount of iterations is able to predict high-fidelity built configurations. Moreover, our technique is language and system agnostic. In essence it only requires monitoring the final result of the build process.

## 9 CONCLUSIONS

We have presented *Reliabilbuild*, an active-learning-based framework to select high-fidelity build configurations using limited number of samples. We implemented two models, one based on *Wisdom of the crowd* and another based on *Bayesian* optimization and evaluate their efficacy in identifying build configurations that are highly likely to build. For the purpose of our evaluation, we collected a large dataset consisting of packages from the E4S package ecosystem. We showed that our models are able to select high fidelity configurations with far fewer samples in comparison to a random exploration. For example, *Reliabilbuild* selects 3× the number of good configurations in comparison to random sampling for several packages including Abyss, Bolt, libnrm, OpenMPI. Our framework is also able to select all the high-fidelity builds in half the number of samples as required by random sampling for packages such as Chai, OpenMPI, py-petsc4py, slepc. We also leveraged our model to provide insights about package incompatibilities, which can be used by both package managers as well as users to identify problematic version pairs and avoid including them. Most importantly, *Reliabilbuild* provides an automatic way to select high-fidelity build configurations with significantly less samples.

## 10 ACKNOWLEDGEMENTS

This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344 (LLNL-CONF-831078). Work at LLNL was funded by the Laboratory Directed Research and Development Program under project tracking code 21-SI-005.

## REFERENCES

- [1] Cargo: The Rust package manager. Online, March 2014. <https://github.com/rust-lang/cargo>.
- [2] PubGrub version solving algorithm implemented in Rust. Online, 2020. <https://github.com/pubgrub-rs/pubgrub>.
- [3] Pietro Abate, Roberto Di Cosmo, Georgios Gousios, and Stefano Zacchiroli. Dependency solving is still hard, but we are getting better at it. In *2020 IEEE 27th International Conference on Software Analysis, Evolution and Reengineering (SANER)*, pages 547–551. IEEE, 2020.
- [4] Pietro Abate, Roberto Di Cosmo, Ralf Treinen, and Stefano Zacchiroli. Dependency solving: a separate concern in component evolution management. *Journal of Systems and Software*, 85(10):2228–2240, 2012.
- [5] Megha Agarwal, Divyansh Singhvi, Preeti Malakar, and Suren Byna. Active learning-based automatic tuning and prediction of parallel i/o performance. In *2019 IEEE/ACM Fourth International Parallel Data Systems Workshop (PDSW)*, pages 20–29. IEEE, 2019.
- [6] Josep Argelich, Daniel Le Berre, Inês Lynce, João P. Marques Silva, and Pascal Rapicault. Solving linux upgradeability problems using boolean optimization. In Inês Lynce and Ralf Treinen, editors, *Proceedings First International Workshop on Logics for Component Configuration, LoCoCo 2010, Edinburgh, UK, 10th July 2010*, volume 29 of *EPTCS*, pages 11–22, 2010.
- [7] Cyrille Artho, Kuniyasu Suzuki, Roberto Di Cosmo, Ralf Treinen, and Stefano Zacchiroli. Why do software packages conflict? In *2012 9th IEEE Working Conference on Mining Software Repositories (MSR)*, pages 141–150. IEEE, 2012.
- [8] James S Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter optimization. In *Advances in neural information processing systems*, pages 2546–2554, 2011.
- [9] Alex Birsan. Dependency confusion: How i hacked into apple, microsoft and dozens of other companies, 2021.
- [10] Christopher Bogart, Christian Kästner, James Herbsleb, and Ferdian Thung. How to break an api: cost negotiation and community values in three software ecosystems. In *Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering*, pages 109–120, 2016.
- [11] Junjie Chen, Ningxin Xu, Peiqi Chen, and Hongyu Zhang. Efficient compiler autotuning via bayesian optimization. In *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*, pages 1198–1209. IEEE, 2021.

- [12] Koen Claessen and John Hughes. Quickcheck: a lightweight tool for random testing of haskell programs. In *Proceedings of the fifth ACM SIGPLAN international conference on Functional programming*, pages 268–279, 2000.
- [13] Bradley E Cossette and Robert J Walker. Seeking the ground truth: a retroactive study on the evolution and migration of software libraries. In *Proceedings of the ACM SIGSOFT 20th International Symposium on the Foundations of Software Engineering*, pages 1–11, 2012.
- [14] Russ Cox. What npm should do today to stop a new colors attack tomorrow, 2022.
- [15] Alexandre Decan and Tom Mens. What do package dependencies tell us about semantic versioning? *IEEE Transactions on Software Engineering*, 2019.
- [16] Alexandre Decan, Tom Mens, and Philippe Grosjean. An empirical comparison of dependency network evolution in seven software packaging ecosystems. *Empirical Software Engineering*, 24(1):381–416, 2019.
- [17] Roberto Di Cosmo. EDOS deliverable WP2-D2.1: Report on Formal Management of Software Dependencies. Technical report, INRIA, May 15 2005. hal-00697463.
- [18] Jens Dietrich, Kamil Jezek, and Premek Brada. What java developers know about compatibility, and why this matters. *Empirical Software Engineering*, 21(3):1371–1396, 2016.
- [19] Jens Dietrich, David Pearce, Jacob Stringer, Amjed Tahir, and Kelly Blincoe. Dependency versioning in the wild. In *2019 IEEE/ACM 16th International Conference on Mining Software Repositories (MSR)*, pages 349–359. IEEE, 2019.
- [20] P. F. Dubois, T. Epperly, and G. Kurfert. Why johnny can't build [portable scientific software]. *Computing in Science Engineering*, 5(5):83–88, 2003.
- [21] Marc Ewing and Erik Troan. RPM Timeline. Online, 1995. <https://rpm.org/timeline.html>.
- [22] Todd Gamblin. Software integration challenges. Technical report, Lawrence Livermore National Lab.(LLNL), Livermore, CA (United States), 2021.
- [23] Todd Gamblin, Matthew LeGendre, Michael R Collette, Gregory L Lee, Adam Moody, Bronis R De Supinski, and Scott Futral. The spack package manager: bringing order to hpc software chaos. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–12, 2015.
- [24] Martin Gebser, Roland Kaminski, and Torsten Schaub. aspcud: A linux package configuration tool based on answer set programming. *Electronic Proceedings in Theoretical Computer Science*, 65:12–25, Aug 2011.
- [25] Martin Gebser, Benjamin Kaufmann, Roland Kaminski, Max Ostrowski, Torsten Schaub, and Marius Schneider. Potassco: The potsdam answer set solving collection. *AI Communications*, 24(2):107–124, 2011.
- [26] Jesus M Gonzalez-Barahona, Gregorio Robles, Martin Michlmayr, Juan José Amor, and Daniel M German. Macro-level software evolution: a case study of a large software compilation. *Empirical Software Engineering*, 14(3):262–285, 2009.
- [27] Jason Gunthorpe. APT User's Guide. Online, 1998. <https://www.debian.org/doc/manuals/apt-guide/>.
- [28] Qiang He, Bo Li, Feifei Chen, John Grundy, Xin Xia, and Yun Yang. Diversified third-party library prediction for mobile app development. *IEEE Transactions on Software Engineering*, 2020.
- [29] Joseph Hejderup and Georgios Gousios. Can we trust tests to automate dependency updates? a case study of java projects. *Journal of Systems and Software*, 183:111097, 2022.
- [30] K. Hoste, J. Timmerman, A. Georges, and S. D. Weirtd. Easybuild: Building software with ease. In *2012 SC Companion: High Performance Computing, Networking Storage and Analysis*, pages 572–582, 2012.
- [31] Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. Sequential model-based optimization for general algorithm configuration. In *International conference on learning and intelligent optimization*, pages 507–523. Springer, 2011.
- [32] Kamil Jezek and Jens Dietrich. Api evolution and compatibility: A data corpus and tool evaluation. *J. Object Technol.*, 16(4):2–1, 2017.
- [33] Donald R Jones. A taxonomy of global optimization methods based on response surfaces. *Journal of global optimization*, 21(4):345–383, 2001.
- [34] Raula Gaikovina Kula, Daniel M German, Ali Ouni, Takashi Ishio, and Katsuro Inoue. Do developers update their library dependencies? *Empirical Software Engineering*, 23(1):384–417, 2018.
- [35] G Kurfert and T Epperly. Software in the DOE: The Hidden Overhead of “The Build”. Technical report, Lawrence Livermore National Laboratory, February 28 2002. UCRL-ID-147343.
- [36] Meir M Lehman. Programs, life cycles, and laws of software evolution. *Proceedings of the IEEE*, 68(9):1060–1076, 1980.
- [37] F. Mancinelli, J. Boender, R. di Cosmo, J. Vouillon, B. Durak, X. Leroy, and R. Treinen. Managing the complexity of large free and open source package-based software distributions. In *21st IEEE/ACM International Conference on Automated Software Engineering (ASE'06)*, pages 199–208, 2006.
- [38] M Douglas McIlroy, J Buxton, Peter Naur, and Brian Randell. Mass-produced software components. In *Proceedings of the 1st international conference on software engineering, Garmisch Pattenkirchen, Germany*, pages 88–98, 1968.
- [39] Wes McKinney et al. pandas: a foundational python library for data analysis and statistics. *Python for High Performance and Scientific Computing*, 14(9):1–9, 2011.
- [40] Harshitha Menon, Abhinav Bhatele, and Todd Gamblin. Auto-tuning parameter choices in hpc applications using bayesian optimization. In *2020 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pages 831–840. IEEE, 2020.
- [41] Claude Michel and Michel Rueher. Handling software upgradeability problems with MILP solvers. In Inès Lynce and Ralf Treinen, editors, *Proceedings First International Workshop on Logics for Component Configuration, LoCoCo 2010, Edinburgh, UK, 10th July 2010*, volume 29 of EPTCS, pages 1–10, 2010.
- [42] Yana Momchilova Mileva, Valentin Dallmeier, Martin Burger, and Andreas Zeller. Mining trends of library usage. In *Proceedings of the joint international and annual ERCIM workshops on Principles of software evolution (IWPSE) and software evolution (Evol) workshops*, pages 57–62, 2009.
- [43] Jonas Mockus. *Bayesian approach to global optimization: theory and applications*, volume 37. Springer Science & Business Media, 2012.
- [44] Phuong T Nguyen, Juri Di Rocco, and Davide Di Ruscio. Mining software repositories to support oss developers: A recommender systems approach. In *IIR*, 2018.
- [45] Phuong T Nguyen, Juri Di Rocco, Davide Di Ruscio, and Massimiliano Di Penta. Crossrec: Supporting software developers by recommending third-party libraries. *Journal of Systems and Software*, 161:110460, 2020.
- [46] Phuong T Nguyen, Juri Di Rocco, Riccardo Rubei, Claudio Di Sipio, and Davide Di Ruscio. Recommending third-party library updates with lstm neural networks. 2021.
- [47] Lina Ochoa, Thomas Degueule, Jean-Rémy Falleri, and Jurgen Vinju. Breaking bad? semantic versioning and impact of breaking changes in maven central. *arXiv preprint arXiv:2110.07889*, 2021.
- [48] Ali Ouni, Raula Gaikovina Kula, Marouane Kessentini, Takashi Ishio, Daniel M German, and Katsuro Inoue. Search-based software library recommendation using multi-objective optimization. *Information and Software Technology*, 83:55–75, 2017.
- [49] Tom Preston-Werner. Semantic versioning 2.0. 0, 2013.
- [50] Python Software Foundation. New pip resolver to roll out this year. Online, March 23 2020. <https://pyfound.blogspot.com/2020/03/new-pip-resolver-to-roll-out-this-year.html>.
- [51] Zhensu Sun, Yan Liu, Ziming Cheng, Chen Yang, and Pengyu Che. Req2lib: A semantic neural model for software library recommendation. In *2020 IEEE 27th International Conference on Software Analysis, Evolution and Reengineering (SANER)*, pages 542–546. IEEE, 2020.
- [52] Chris Tucker, David Shuffleton, Ranjit Jhala, and Sorin Lerner. OPIUM: Optimal package install/uninstall manager. In *International Conference on Software Engineering (ICSE)*, 2007.
- [53] Andre Van Der Hoek, Richard S Hall, Dennis Heimbigner, and Alexander L Wolf. Software release management. *ACM SIGSOFT Software Engineering Notes*, 22(6):159–175, 1997.
- [54] Julien Villemontex, Emmanuel Vazquez, and Eric Walter. An informational approach to the global optimization of expensive-to-evaluate functions. *Journal of Global Optimization*, 44(4):509–534, 2009.
- [55] Natalie Weizenbaum. PubGrub: Next-Generation Version Solving. <https://medium.com/@nex3/pubgrub-2fb6470504f>, April 2 2018.
- [56] Shengzhe Xu, Ziqi Dong, and Na Meng. Meditor: inference and application of api migration edits. In *2019 IEEE/ACM 27th International Conference on Program Comprehension (ICPC)*, pages 335–346. IEEE, 2019.