Graph-Based Orchestration of Heterogeneous AI Models: A Control Node Approach to Composite Intelligence

G. Michael Youngblood, Filip Dvořák, Slavomir Svancar, Tomáš Balyo, Michal Ficek and Martin Dousek

Filuta AI, Inc., 1606 Headway Cir STE 9145, Austin, TX 78754, USA

{michael, filip, slavo, tomas, michael, martin}@filuta.ai

Abstract

We present a graph-based framework for Com-1 posite AI systems that integrate multiple AI tech-2 niques to address complex problems. We describe 3 an architecture using specialized control nodes 4 to orchestrate diverse models including planning, 5 machine learning, constraint programming, and 6 large language models. The system employs di-7 rected acyclic graphs and behavior trees to manage 8 model relationships and execution sequencing. Our 9 methodology encompasses data connectors, com-10 putational services, and a containerized deploy-11 ment framework supporting hierarchical, sequen-12 tial, and concurrent model composition patterns. 13 Through case studies, we demonstrate how this ar-14 chitecture enables the creation of composite sys-15 tems that leverage the strengths of multiple AI tech-16 niques within a unified operational schema, over-17 coming the limitations of single-model approaches. 18

19 1 Introduction

There are many techniques for modeling real-world prob-20 lems, and no single technique will dominate as the most 21 efficient across most real-world problems for the near fu-22 ture. Instead, modeling real-world problems can be seen as a 23 composition of chained, concurrent (parallel or interleaved), 24 and/or nested techniques/models, using the most efficient 25 technique/model for each of the sub-problems into which 26 the real-world problem decomposes. Significant validation 27 of this composition approach to solving industrial problems 28 can be found in the studies performed by Gartner [Gartner, 29 2022], where Composite AI has emerged as one of the most 30 promising techniques, following the under-performance of 31 deep learning and machine learning in applied projects. 32

The system and methods presented here provide a common 33 ground for the composition of various models and techniques, 34 which we will just refer to simply as just "models" going for-35 ward, into a system. These models include AI Planning, ma-36 chine learning, scheduling, graph theory, constraint program-37 ming, linear programming, large language models (LLMs) 38 such as GPT [Floridi and Chiriatti, 2020], multimodal foun-39 dational models, and any other type of element that using 40

computation can perform a distinct set of functions producing output from input. 41

This work is motivated by creating systems composed of 43 many models to address the challenges single models face 44 in being able to perform many disparate functions outside of 45 their design/training. Models cannot reason about the opti-46 mal plan of actions to achieve a goal, dispatch and monitor 47 action execution, act on state information, diagnose failures, 48 and adapt to unexpected events all in a single model. Current 49 automation challenges also involve connecting all actors in 50 the system and collecting their real-time operational data. Ef-51 ficient and reliable control of a target system to achieve goals 52 requires reasoning about this data using an appropriate com-53 position of models leading to appropriate action decisions. 54

We introduce a methodology including key building blocks 55 for constructing automation systems designed as a graph with 56 the flexibility to transform into many graph forms for deploy-57 ment. The key architectural invention is the introduction of a 58 data and model aware, logic-driven control node central to 59 defining a composition that unifies a plurality of disparate 60 models. These compositions can be used to generate code, 61 be containerized, and then deployed autonomously. 62

The Composite AI system provides capabilities to incorporate individual models, allowing the human user to augment the models, compose the models into the operational schema with well-defined inputs, outputs, and goals, and generate stand-alone intelligent automation services that act towards user-defined goals based on the model composition and observations of the world.

2 Graphs

Our Composite AI system's foundation is a graph structure, which models pairwise relationships between entities through vertices (nodes) and edges (links). This representation offers significant implementation advantages, including clear visualization of relationships, efficient algorithmic processing, modular design, and adaptability to dynamic changes. Graphs effectively capture complex hierarchical structures, aid in problem decomposition, and support various data storage formats, making systems more efficient and scalable.

2.1 Directed Acyclic Graph (DAG)

A Directed Acyclic Graph (DAG) features directed edges 81 with specific orientation and contains no cycles. This struc-

70 71

72 73 74

79

80

ture excels at managing dependencies and order, ensuring
tasks execute only after prerequisites are met. The absence
of cycles eliminates circular dependencies, simplifying resolution and preventing deadlocks. DAGs support efficient optimization algorithms and provide clarity in representing complex systems, making them ideal for system management and
control.

90 2.2 Behavior Trees (BT)

Behavior trees are hierarchical graphs that define autonomous
agent behavior in a tree-like structure. They consist of three
node types:

- Gomposite nodes (→ and ?) control execution flow by organizing other nodes
- Decorator nodes modify child node behavior
- Leaf nodes perform actions or checks

Behavior trees enable hierarchical problem decomposition, 98 allowing higher-level trees to incorporate lower-level trees 99 without requiring detailed knowledge of subproblems. They 100 can integrate planning and diagnosis capabilities, with each 101 102 action represented as a node whose evaluation reflects the ac-103 tion's execution status. Within the framework, behavior trees function as deterministic policies that observe the world and 104 initiate appropriate actions. Figure 1 illustrates an example 105 behavior for monitoring and adapting a plan in execution. In 106 the example, blue nodes have the following behavior: 107

- 108 \rightarrow (and) evaluates to:
- success if all its children evaluate to success,
- failure if at least one node evaluates to failure, and
- running if at least one node evaluates to running
 while none evaluates to failure.
- ? (or) evaluates to:

114

115

- success if at least one child evaluates to success,
- failures if all children evaluate to failures, and
- running if at least one child evaluates to running and none evaluates to success.



Figure 1: Example Behavior Tree

The key advantages of behavior trees include verifiability,
human-readability, and visualization capabilities, providing a
common platform for sharing architectural concepts and domain knowledge among stakeholders.

Within a framework, a behavior tree can be seen as a (deterministic) policy that observes the world with each tick of the tree and may send actions, such as "Plan." Planning can be a subproblem itself, implemented as a policy through a solver that takes state, actions, and a goal to develop a plan for execution.

Behavior trees are verifiable, human-readable, and easy to visualize. They provide a common ground for sharing ideas, architectures, and domain knowledge among domain experts, product managers, and AI systems.

3 Methodology

Learning a single end-to-end model is rarely computationally feasible. Instead, end-to-end models are typically composed of multiple heterogeneous models, some of which are machine-learned while others are human-built using AI and operations research modeling techniques. The system and methods described in this paper support the design, development, and deployment of these composite AI systems. 133

The system goal is the automated assembly of data con-140 nectors, domain and goal-based problem models, computa-141 tional services (e.g., solvers, reasoners, large-language mod-142 els, etc.), and appropriately configured control nodes (logic-143 gated demultiplexers, multiplexers, and pass-through using 144 data and model knowledge) using a graph-based organization 145 into a service supported by an externally interfacing API that 146 can be plugged into a target system for automation as a Com-147 posite AI system. The Composition Framework shown in 148 Figure 2 illustrates an abstract graph of the principal compo-149 nents in making compositions (the blue parts). Online and/or 150 offline data is brought into the system and adapted for dis-151 tribution to one or more models organized into one or more 152 stages. Novel control nodes can flexibly be configured to dis-153 tribute data into and out of model in to, between, and out 154 of stages. As shown, some models require computational 155 pipelining. At the end of all stages the composition produces 156 output. This composite core is wrapped in a service API and 157 designed for deployment in a container or set of containers to 158 provide the configured service both in a testing and produc-159 tion environment. 160

3.1 Online and Offline Data Connectors

When building real-time intelligent services, it is necessary to process both offline and online data. Offline data, which are typically used for model learning, can be quite large and are not significantly constrained by processing time. In contrast, online data are expected to arrive in real-time when the service is deployed, and the results of processing the real-time data are going to be produced and published by the service. 163

While the connectors are critical components of the plat-169 form and composed services, the implementation effort has 170 already been invested in by several open-source data con-171 nectors, such as Airbyte [Airbyte, 2023]. The system builds 172 upon an existing framework of connectors developed by oth-173 ers and extends their functionality by adding preprocessing 174 data transformations relevant to learning and/or deploying the 175 models. 176

161



Figure 2: Composition framework. Each control node is unique being configured for their specific downstream components through introspection or manually, but use a common code base.

177 3.2 Models

The core workers in a composite AI system are the mod-178 els. Models may be pre-trained with data, augmented with 179 new data, learn from zero using new data, or only algorith-180 mic in nature (i.e., a technique). When models need to be 181 augmented or developed from the ground up then the model-182 specific techniques apply such as Automated Domain-driven 183 Model Synthesis [Balyo *et al.*, 2024]. By applying various 184 filters and transformations to a dataset ingested from a data 185 connector, the system (optionally aided by a user) can define 186 the state space or other features, which are then used to learn 187 the model. The same dataset can be used to define state spaces 188 or inputs for multiple different modeling approaches. For ex-189 ample, we can model a predictive task as a dependency of the 190 engine temperature on speed, road incline, and total weight, 191 while simultaneously modeling vehicle logistics using classi-192 cal planning [Ghallab et al., 2004]. 193

194 3.3 Computational Services

Some models require computational support to execute. For 195 example, a planning model needs to be run through a solver 196 with a problem specification. A GPT (Generative Pre-trained 197 198 Transformer) model needs to have its weights loaded into the configured transformer architecture and be run with a prob-199 lem prompt. One case study of the system supports model-200 ing with classical planning using machine learning [Dvorak 201 et al., 2021], as well as constraint modeling using CP-SAT 202 in [Google, 2022] and machine learning [LeDell and Poirier, 203 2020]. 204

3.4 Control Nodes: Logic-Gated Multiplexers, Demultiplexers, and Pass-throughs

The key to creating complex composite AI systems is providing a control mechanism to combine, select, transform, deconstruct, and configure how the data flows into and out of the models in each stage until final output. Each control node is configured to marshal the data into or out of a set of models and configure it as required using the following configuration elements. 213

- 1. Permeability: Defines how data flows through the gate. 214
 - (a) Fully permeable: Data is a complete pass-through to one or more models downstream. 216
 - (b) Partially permeable: Data is a partial pass-through to one or more models downstream and may be combined, selected, or transformed before being passed to the model(s).
 - (c) Impermeable: Data is not passed-through and will
 be combined, selected, or transformed before being
 passed to one or more models downstream if at all.
- 2. **Combination**: Define how inputs are combined before being passed to one or more models downstream.
- 3. **Selection**: Define how selected inputs are passed to one or more models downstream if at all (e.g., filtered). 227
- 4. **Transformation**: Define how inputs are altered before 222 being passed to one or more models downstream. 222
- 5. Downstream Model Configuration: Define specific model configurations for each model as needed for each input. For example, each input into a LLM may specify some specific hyper-parameters associated for processing the input.
- 6. **Orchestration**: Defines conditions of received outputs 235 and when/how to proceed in processing. Orchestra-236 tion may only require a few responses instead of all re-237 sponses from a set of upstream models before proceed-238 ing. It may require specific properties to the responses 239 before proceeding. However, it should never perma-240 nently halt processing and any resolutions needed to pro-241 ceed should be able to be kicked off by the gate. 242

In a typical case study as shown in Figure 3, the control 243 node is implemented in a computer programming language 244

into four areas of functionality. Data comes into a Data 245 Combination, Selection, and Transformation Logic compo-246 nent that applies logic rules to the incoming data to combine, 247 select (filter in or out), and transform the data as specified. 248 A Model Configuration and Orchestration Logic component 249 receives information about the models it can use and is con-250 figured through logic rules on how to set up, execute, and 251 orchestrate the downstream models. These two components 252 may share information that informs logic rules for processing. 253 The Data to Model component provides the data to the ap-254 255 propriate downstream component(s), which may be the system output. The Model Control component configures the 256 downstream models, if they exist, and kicks off their opera-257 tion as defined by the logic rules. Control Nodes are aware 258 through introspective or manual configuration of their down-259 stream components and parameters. Control nodes always 260 have incoming data and should have some outgoing data but 261 may not have any model input or control. Any logical speci-262 fication system can be used to implement a control node. 263



Figure 3: Control Node Architecture

One example of how control nodes are used in a stage is 264 to reduce LLM hallucinations, generated incorrect or fabri-265 cated information, by passing the same prompt into each of 266 a set of LLM models (e.g., ChatGPT, Llama, and Gemini) in 267 parallel. The outputs are then transformed into an intersect-268 ing set so that only outputs that are in fuzzy match are passed 269 to the next stage. Information that is not the same in each 270 response is filtered reducing inconsistent aspects, which are 271 often manifested as LLM hallucinations. 272

Control nodes can output to other control nodes and be 273 nested in model flows as shown in Figure 1. Control nodes fit 274 neatly into the system graph as nodes before and after models 275 and model processing nodes. Most times system control is 276 best defined as a DAG or BT, but in some case studies a di-277 rected graph may be used with cycles. These cycles typically 278 279 specify a loop between control nodes so that information may be processed multiple times by the same set of models. 280

Compositions may be designed to operate under various modalities. One common set is the train and test (i.e., learn and perform, explore and exploit). Control nodes can be specified with different logic configurations to support each desired mode of composition usage (e.g. using environmental variables).

4 Architectural Solution Composition

Architectural Solution Composition in practice involves 288 building a graph (behavior tree in this example) that combines multiple models with corresponding solvers, goals, and 290 APIs into a service. 291

Although the models share a single observation space, each 292 model only considers its own local state space. Each model 293 can be tied to a computation service (e.g., solver), which may 294 continuously work on problem instances generated from the 295 observations whenever it is triggered. The system can use 296 any type of software defined computation services including 297 planning solvers in the Unified Planning Framework [Micheli 298 et al., 2025] and the wide range of AI techniques. 299

We often encapsulate reasoning and action execution into the nodes of a behavior tree providing the grounds to define a composed automated reasoning, planning, and execution system in terms of tree composition. The tree composition is an operationalization of multiple models together with the corresponding solvers, actors, goals, and APIs into a service. Any two or more models in a composition are related either by

- 1. nesting (hierarchical), when one model is fully nested within another, and the subsystem operating with the outer model can invoke the subsystem operating the inner model,
 307

 309
 309

 301
 309

 302
 310
- precedence or chaining (sequential), when a subsystem
 for the earlier model can be fully executed before running the subsystem for the later model, or
- concurrent (interleaving and/or parallel), when subsystems for both models can run concurrently and interfere with each other, for example by exchanging their best solutions.
 314
 315
 316
 317

These three model composition approaches are available 318 to the solution architect through the control nodes, and it is 319 upon the architect to design a composition appropriate for the 320 given real-world problem. The models share a single obser-321 vation space, the control nodes picking only the data relevant 322 for them (e.g., in a food delivery system example, vehicle 323 routing with time windows precedes running a planner). We 324 first produce optimal assignments of orders to drivers, which 325 is equivalent to adding control knowledge to the more general 326 planning problem, which is solved consequently. We consider 327 each model to be tied with a computational service subsystem 328 as appropriate (e.g., plan with a planner, problem solve with a 329 solver), which keeps solving the problem instances generated 330 from the observations every time it is triggered. 331

The concept of model composition in the system is another tool typically provided to the solution architect user. If desired, the user can build and compose models to guarantee optimal solutions. There can also be cases when optimality is infeasible or not desired, and the goal is to find solutions quickly with a reasonable distance from optimality. 332

Service code can be generated by introspectively filling out template code of various components of the composition, or by configuring pre-build services (e.g., using JSON configuration files, environment variables, or API calls). In implementation there typically exists a code/service template for every type of composition node. The solution is composed of 343

a main orchestration service, that contains knowledge about 344 all the components in the composition. The individual nodes 345 are either directly part of the orchestration service, in case 346 of simple nodes, or the orchestration service has lightweight 347 proxy wrappers for a dedicated microservice that are logically 348 part of the composition. The orchestration service communi-349 cates with the microservice using this proxy and can either 350 directly fetch and push data to it via this proxy or can use it to 351 configure the wrapped service to specify data inputs and out-352 puts. Individual nodes are decoupled, language-independent, 353 and are flexible to be able to be part of a wide range of com-354 positions. 355

All the parts of the solution templates, including the or-356 chestration service and all the engineered packages should 357 be packed into container technology (i.e., Docker) by design. 358 After code generation, all that is needed is to execute the al-359 ready created build commands to create the containers and 360 service manifest. These containers then can be easily ver-361 sioned, duplicated, and deployed either locally or to a cloud 362 solution, and can be used in various environments, such as 363 testing, staging, demo, or production. 364

The Filuta AI system implements all of these aspects into a commercial platform used to provide solutions to real-world problems such as video game quality assurance evaluations.

368 5 Deployment

Deploying services produced by the system is a push-button action for the user. It involves a cloud-agnostic or computational platform deployment of containerized microservices at the backend. Once deployed, users can monitor the deployment status and interact with it at the defined network address through the API defined in the composition.

Final deployment also contains implicit and explicit tools 375 and services not directly part of the data pipeline. These ser-376 vices allow us to incorporate functional and non-functional 377 requirements to the solution, such as High Availability, ex-378 tensive monitoring and logging capacities, authorization, and 379 authentication mechanisms. Tools such as monitoring are im-380 plicitly part of every deployment and are not visualized in the 381 Composition UI representation. An example of explicit sup-382 383 porting tools is the High Availability Service that will deploy the Composition multiple times in multiple regions to serve 384 as a failback if one instance fails, to allow faster and more 385 scalable usage for customers from different regions, or for 386 many customers in one place. 387

- 388 Deployed services can be used in two separate ways:
- Dedicated for a composition to which they belong
- Distributed across multiple service environments, customer solutions, or even multiple customers.

Dedicated services are, for example, models dependent on 392 the environment where they are used or if the customer wants 393 dedicated services running for them. Distributed services can 394 be either shared data sources that are the same for all the cus-395 tomers, for example a weather service, or resource-heavy ser-396 vices that would be cost prohibitive to use for a single deploy-397 ment, such as an LLM service. Customer isolation is always 398 399 guaranteed by the solution either by using dedicated services or by not sharing context in the distributed services. 400

Composition is the core part of a complete system as shown in Figure 4 with four main components: 402

1. Model Editor allows the expert user to define data inputs, 403 and how they are transformed into an ingestible format. For 404 example, tabular data for predictors/ML models, or time se-405 ries for symbolic models. It then provides a way to plug those 406 connectors/transformers into an ML or symbolic modeling 407 approach. It is possible for the expert user to interactively 408 adjust and augment the proposed model until they are satis-409 fied with its quality. This could mean accuracy for ML, or 410 human-check and test validation against input data for sym-411 bolic models. 412

Models are used in Compositions and have the following413properties: 1) Defined input specifications, 2) Defined output414specifications, 3) Defined model configuration parameters, 4)415Defined model controls, 5) Defined computation and memory requirements, and 6) Characterization of performance on416data.418

2. Composition Editor takes multiple models and offers di-419 verse options to combine them using control nodes to imple-420 ment a graph that can be executed as a DAG, behavior tree, 421 or other forms as well as using data connectors to provide 422 shared observation space. The editor encodes both the de-423 composition of the mathematical structure representing the 424 real-world problem and the operationalization of the model 425 with the behavior of a general agent. This feature allows for 426 quick iteration, learning, and versioning of agents that pro-427 vide intelligent automation services. The editor also specifies 428 I/O 429

The building blocks of a composition are explicitly ver-430 sioned to support breaking changes-if one service changes 431 their API and other dependent service needs to adjust, the 432 system can keep one composition that contains services that 433 have not adjusted yet on an old version and another composi-434 tion that already has all the changes incorporated at a newer 435 version. This also supports falling back to older versions if 436 needed, rolling out features gradually, and other cases where 437 explicit versioning supports better system management. 438

The building blocks are used in case studyies to 1) con-439 struct the behavior tree or other executable graph in Python 440 (or another OOP language based on customer needs, such 441 as C++ for game testing integrations) and 2) list the deploy-442 ment requirements of the components necessary to run the 443 composition. The behavior tree represented by the composi-444 tion in a case study has shown to be a low-latency (1-5ms), 445 minimal-computation tree structure that constantly runs when 446 deployed and manages the other components (containers) 447 considered to be serviced in this interaction. 448

3. Simulation Testing Manager deploys compositions by
mapping them into a simulation representative of the target
system. Control of the simulation can run the composition
through testing phases with the simulated system under any
type of definable conditions.449
450

4. **Deployment Manager** deploys compositions by mapping them 1:*n* into running instances, each with its own configuration and connected to a target system. For example, we can assign a location (e.g., DNS name or IP number) and an authentication token as minimal configurations. The deployment manager continually monitors, terminates, and pro-



Figure 4: An end-to-end system for Specification, Codification, and Deployment of Composite AI Systems

vides other features expected from cloud-deployed serviceswith minimal downtime.

Composite AI systems are designed to be systemagnostic—they function the same way regardless of whether they're connected to the actual target system or a simulation of it. For testing purposes, users can connect their compositions to existing simulators (if available) or use configurable generic simulators that mimic the target system's behavior.

468 6 Target System Application

We describe an approach for building a general hierarchical
autonomous control architecture using a composition. The
approach combines well-known techniques from satisfiability
theory, automated planning, automated diagnosis, constraint
programming, reinforcement learning, and behavior trees.

A user provides a **target system** that accepts actions and provides a perception of it and its world's state. The user works with the backend system through a **user interface** to configure one or more domain models and a composite system of AI-based components organized by a graph (e.g., behavior tree) running through tests and deployment to solve problems through action-space control.

481 As a running example, we will describe the components of a specific case study of this system that has been set up to 482 solve the Towers of Hanoi game as the target system. The 483 Towers of Hanoi is a puzzle game that involves moving discs 484 from one place to another. Imagine you have n rods, where 485 n is typically three, and a stack of different-sized discs. The 486 goal is to move all the discs from one rod to another, fol-487 lowing a few rules. First, you can only move one disc at a 488 time. Second, you can never put a bigger disc on the top of 489 a smaller one. And third, you can use the other rods to help 490 you move the discs around, but you can only place a disc on 491 top of another disc or on an empty rod. 492

The **deployed composition** contains two complementary models: a symbolically trained Towers of Hanoi Model that captures domain knowledge in PDDL format, and a Reinforcement Learning (RL) Move Evaluator trained on historial gameplay data to predict move quality scores. The models were learned by taking raw logs from a system, an expert's domain knowledge, existing domain knowledge (e.g., existing PDDL descriptions), and other sources that provide do-500 main knowledge through capture collecting into a knowledge 501 base. Knowledge capture primarily involves preparing a set 502 of predicates that describe the domain under interest follow-503 ing any method of knowledge capture from a consultant or 504 domain expert. A predicate is a function of a set of parame-505 ters that returns a Boolean as an answer. The knowledge base 506 may also contain a description of predicates and their mean-507 ing extracted from the target system's raw logs. Predicates 508 apply to a specific type of object, or to all objects. Predicates 509 are either true or false at any point in a plan and when not 510 declared, they are assumed to be false. In our running exam-511 ple, the knowledge base may contain information such as the 512 definition of useful predicates: *discONdisc(disc1, disc2)* 513 means disc1 is on disc2, smaller(disc1, disc2) means disc1 514 is smaller than disc2, discClear(disc1) means disc1 is on 515 the top and nothing lies on it. It may also define the goal to 516 be achieved, such as all discs lie on the last rod, composed of 517 predicates that are all true (e.g., discONrod(disc1, rod3), 518 discONrod(disc2, rod3), and discONrod(disc3, rod3)). 519 It may also define the initial state of the problem using predi-520 cates that describe the input scenario and all combinations of 521 disc and their sizes and location. 522

Information from the target system flows into and out of 523 the **deployed composition** through the data connectors. This 524 takes perception data from the external world state of the sys-525 tem or other systems that may accurately report data from 526 the environment of the target system. This could be any ex-527 ternal information relevant to the domain of interest, such as 528 weather, traffic conditions, general knowledge from the outer 529 world, etc. The target system may provide telemetry or op-530 erational logs, such as a raw stream/batch of data collected 531 from system components. For our running example, keys 532 pressed in the game by a human player. The composition 533 output data connector transmits action information to the sys-534 tem-transformed planned actions that lead to actual system 535 change. In the planned action for our running example, an ac-536 tion can be *MoveDiscFromRodToRod*, which translated 537 to the target system can mean "press arrow up, press arrow 538 left, press arrow left, press arrow left, press arrow down" as a 539 sequence of virtual keyboard key presses. 540

A data connector takes in data in some form, digests 541

the data from the source (e.g., document, Google Drive, 542 SFTP, stream, websocket, Parquet, AWS), and passes it in 543 a defined structured data format to be processed by other 544 entities. In our running example, this could involve using 545 Keboola or some data platform to interchange data from the 546 user's game into a format defined by the predicate store. The 547 configuration may include format specifications, connectors, 548 logins, secrets, etc. A data connector takes in predicate 549 definition raw data, and transforms it given the predicate 550 definition creating predicates from the input log stream. In 551 our running example, Rod1 = [disc3, disc2] is converted 552 discONrod(disc3, rod1), discONrod(disc2, rod1), 553 to discONdisc(disc2, disc3), smaller(disc2, disc3),554 discClear(disc2).555

Within the deployed composition, the PDDL model pro-556 vides guaranteed valid moves and optimal planning capabili-557 ties, while the RL model contributes learned heuristics about 558 move efficiency and common solution patterns. The input 559 data connector passes state information into a partially per-560 meable control node that implements sophisticated orchestra-561 tion logic. This control node transforms the raw game state 562 into two formats (PDDL predicates for the planner and fea-563 ture vectors for the RL model), orchestrates parallel execution 564 565 of both models, combines outputs by using the planner's valid moves as candidates and the RL model's scores to rank them, 566 and selects the highest-scored valid move as output, falling 567 back to the planner's first suggestion if RL scoring fails. 568

A data connector takes information from the composition 569 output and converts it to an action description with param-570 eters, which describes actions to be executed in the target 571 system to achieve the desired goal. In our running example, 572 something like move(d1, d2, d3) means to move disc1 from 573 disc2 to disc3. The output data connector interfaces with the 574 target platform to execute all action commands on that plat-575 form, thus affecting the state through action. 576

The **testing and deployment user-interface** is the frontend, interactive tool that works with the user. The interface again consists of four main parts as illustrated in Figure 4:

 Model development using a text-based language augmented by a visual programming language. In one case study, the modeling language is PDDL text and Blockly provides a visual builder. The user iteratively refines the model through testing and deployment in a digital simulation of the target system and/or the target system. Model debugging tools also provide insight and feedback for iterative refinement.

Models may be completely synthesized from raw observa-587 tional data using a method such as the Automated Domain-588 driven Model Synthesis [Balyo et al., 2024]. The interface 589 tools can then be used to manipulate, edit, add, and remove 590 actions within the synthesized model. The goal is a model 591 containing information from domain synthesis on which ac-592 tion to use or not to use for control with the target system. 593 User augmentation may refine the model and improve under-594 595 standability by renaming predicates, goals, and actions.

2. *Composition development* using a text-based language augmented by a visual programming language. In one case study, the composition language is text describing a behavior tree and Blockly provides a visual builder. A behavior tree visualizer also provides additional insights. The user iteratively refines the composition and model(s) through testing and deployment in a digital simulation of the target system and/or the target system. Composition debugging tools also provide insight and feedback for iterative refinement.

The goal is given one or more models, data connectors, and control nodes, structure the composition workflow in a way that all components cooperate and achieve the desired system performance. 608

In our example, the composition obtains the puzzle state, then orchestrates both the PDDL planner and RL evaluator through a control node to generate optimally-ranked move sequences, which are executed on the target game.

3. *Simulation testing* provides pre-deployment feedback on the composition and underlying models in action. 614

4. *Deployment* provides the ability to push control by the composition to the target system. This includes live performance monitoring of the system in action and supports defining and observing performance through evaluation metrics.

The **Deployed Model Runtime** takes the deployment 619 setup, description of container images, image registry, rele-620 vant services, external endpoints, API's, data locations, and 621 the input stream of target system predicates to deploy the 622 model on a cloud-based platform. It can also run on any com-623 pute platform. It will run services and let them communicate 624 according to description given by composition and ingest in-625 put data (pre-processed predicates) and output data (actions 626 that lead to goal). The system output is a stream of actions 627 to be executed on the target system. It is up to the target 628 system to make the actions happen in the real world given 629 action input. In our running example, this system could be 630 a low-level Azure cloud deployment configurations of indi-631 vidual composition components (planner, data store, serving 632 APIs), networking and interoperability settings, VPNs, con-633 tainer registry settings, endpoints definition, etc. needed to 634 startup and execute a composed AI system to play the Towers 635 of Hanoi. 636

7 Conclusion

In conclusion, this paper presents a comprehensive frame-638 work for Composite AI systems that effectively integrates 639 multiple AI techniques to solve complex real-world prob-640 lems. The graph-based architecture, featuring specialized 641 control nodes, enables the orchestration of diverse models 642 including planning, machine learning, constraint program-643 ming, and large language models within a unified operational 644 schema. By supporting hierarchical, sequential, and con-645 current model composition patterns, the system overcomes 646 the limitations of single-model approaches while providing 647 flexible deployment options through containerization. The 648 methodology encompasses data connectors, computational 649 services, and control mechanisms that together form a power-650 ful platform for designing, testing, and deploying intelligent 651 automation services. As demonstrated through case studies 652 like the Towers of Hanoi example, this approach enables the 653 creation of robust composite systems that can observe, rea-654 son about, and act upon complex environments, representing 655 a significant advancement in applied AI for industrial appli-656 cations. 657

637

658 References

- [Airbyte, 2023] Airbyte. Reflecting on 2023 (and what's in store for 2024), 2023.
- [Balyo *et al.*, 2024] Tomáš Balyo, Martin Suda, Lukáš
 Chrpa, Dominik Šafránek, Stephan Gocht, Filip Dvořák,
- Roman Barták, and G Michael Youngblood. Planning do main model acquisition from state traces without action
- 665 parameters. *arXiv preprint arXiv:2402.10726*, 2024.
- [Dvorak *et al.*, 2021] Filip Dvorak, Anil Agarwal, and Niko lay Baklanov. Visual planning domain design for pddl us-
- ing blockly. In *ICAPS 2021 Demonstrations*, Sugar Land,
 TX 77478, United States, 2021. Schlumberger.
- ⁶⁷⁰ [Floridi and Chiriatti, 2020] Luciano Floridi and Massimo
 ⁶⁷¹ Chiriatti. Gpt-3: Its nature, scope, limits, and conse⁶⁷² quences. *Minds and Machines*, 30:681–694, 2020.
- ⁶⁷³ [Gartner, 2022] Gartner. What's new in artificial intelligence⁶⁷⁴ from the 2022 gartner hype cycle, 2022.
- ⁶⁷⁵ [Ghallab *et al.*, 2004] Malik Ghallab, Dana Nau, and Paolo
 ⁶⁷⁶ Traverso. *Automated Planning: Theory and Practice*. The
 ⁶⁷⁷ Morgan Kaufmann Series in Artificial Intelligence. Mor-
- 677 Morgan Kaufmann Series in Artificial Inte 678 gan Kaufmann, Amsterdam, 2004.
- ⁶⁷⁹ [Google, 2022] Google. Or-tools google optimization
 ⁶⁸⁰ tools, 2022. Google's software suite for combinatorial op ⁶⁸¹ timization.
- [LeDell and Poirier, 2020] Erin LeDell and S. Poirier. H20
 automl: Scalable automatic machine learning. In *Pro- ceedings of the AutoML Workshop at ICML*, volume 2020, 2020.
- [Micheli et al., 2025] Andrea Micheli, Arthur Bit-Monnot, 686 Gabriele Röger, Enrico Scala, Alessandro Valentini, Luca 687 Framba, Alberto Rovetta, Alessandro Trapasso, Luigi 688 Bonassi, Alfonso Emilio Gerevini, Luca Iocchi, Felix In-689 690 grand, Uwe Köckemann, Fabio Patrizi, Alessandro Saetti, Ivan Serina, and Sebastian Stock. Unified planning: Mod-691 eling, manipulating and solving ai planning problems in 692 python. SoftwareX, 29:102012, 2025. 693