

Deep reinforcement learning for process design: Review and perspective

Qinghe Gao¹, Artur M. Schweidtmann^{1,*}

¹ Delft University of Technology, Department of Chemical Engineering, Van der Maasweg 9, Delft 2629 HZ, The Netherlands

Abstract: The transformation towards renewable energy and feedstock supply in the chemical industry requires new conceptual process design approaches. Recently, deep reinforcement learning, a subclass of machine learning, has shown the potential to solve complex decision-making problems and aid sustainable process design. However, its suitability in static process design still needs to be examined. We discuss the advantages and disadvantages of reinforcement learning for process design. Then, we survey state-of-the-art research through three major elements: (i) information representation, (ii) agent architecture, and (iii) environment and reward. Moreover, we discuss perspectives on underlying challenges and promising future works to unfold the full potential of reinforcement learning for process design in chemical engineering.

Keywords: Process synthesis, deep reinforcement learning, machine learning, chemical engineering, artificial intelligence, graph neural network

1 Introduction

The chemical industry is facing a rapid paradigm shift towards a circular economy based on renewable energy and feedstock supply [1, 2]. This poses several challenges for conceptual process design due to the increasing complexity of the design task, the lack of experienced engineers, and the pressure on improving sustainability and profitability while shortening development times. Thus, there is a need for new methodologies and tools that support engineers to design sustainable processes in a more efficient way.

Computer-aid process design (CAPD) is widely used in process systems engineering (PSE) for conceptual process design [3, 4], which can be classified into three methodologies as illustrated in Figure 1: (i) heuristic-based methods; (ii) optimization-based methods; and (iii) the emerging field of generative artificial intelligence (AI). Heuristic-based methods rely on a set of rules derived from experience, insights, and engineering knowledge, making them the most commonly used approaches due to their ease of application. Optimization-based approaches are commonly used to identify optimal design and operating variables for given process structures. Particularly, derivative-based optimization algorithms are already available in commercial process simulation software to determine those. To determine optimal process structures, superstructure-based methods are the de facto state of the art, where possible process alternatives are modeled and subsequently solved by mixed-integer nonlinear optimization (MINLP)

Corresponding author: *A. M. Schweidtmann, E-mail: a.schweidtmann@tudelft.nl

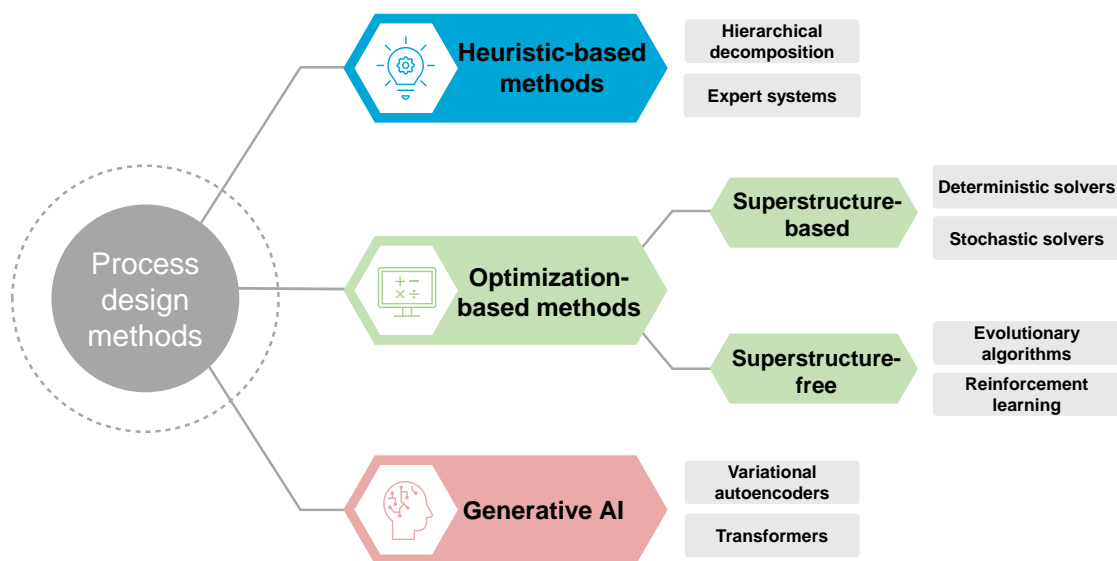


Figure 1: Overview of computer-aided process design methods.

methods [5, 6]. Within this paradigm, two solver types are typically utilized: deterministic (e.g., branch and bound algorithms like BARON or MAiNGO), and stochastic solvers (e.g., genetic algorithms). While superstructure methods have been very successful in PSE, they also have many shortcomings that limit industrial applications [7], including (manual) setup of all process alternatives, the need to implement process models in an optimization environment, and the difficulties of solving resulting MINLPs. Superstructure-free methods dispense the need for predefined superstructures. For example, evolutionary algorithms, typically adopt a two-level decomposition approach [6]. First, they generate alternative flowsheets which are then evaluated through an optimization algorithm. Finally, a few very recent works proposed the use of generative AI methods for the generation of process structures [8, 9]. However, those works require large training data and are not the focus of this review.

Recently, deep reinforcement learning (RL) has shown its potential to solve complex sequential dynamic decision-making problems at human-like or even superhuman level [10, 11, 12]. RL is a computational approach for goal-directed learning and decision-making through the direct interaction of an agent with its environment [13]. RL is primarily developed to solve discrete-time dynamic optimization problems formulated as Markov Decision Processes. Consequently, RL is based on the Bellman optimality equation, which is similar to the Hamilton-Jacobi-Bellman (HJB) equation and Pontryagin’s Maximum Principle (PMP) for continuous state and action spaces in the control theory. Also, RL has seen its first applications in chemical engineering for process control [14, 15] and scheduling [16, 17]. Notably, Yokogawa is even using RL to operate an industrial chemical process since 2022.

This review and perspective paper aims to provide a critical examination of the application of RL

for process design. In the context of process design, RL can be considered a superstructure-free and model-free method, which iteratively places unit operations with corresponding design and operating variables. It evaluates the resulting flowsheets at every iteration and aims to maximize the objective functions. In the recent literature, there have been a few first steps towards applying RL as a static optimization algorithm for stationary process design including absorption–stripping process [18], energy systems design [19, 20], unit operation design [21], separation process [22, 23, 24, 25, 26, 27, 28, 29], solvent extraction process design [30], single mixed-refrigerant process design [31], and synthesis reaction process design [32, 33, 34].

The use of RL for stationary process design is controversial in scientific discussions. Most notably, RL is better suited for dynamic, sequential-decision problems rather than static ones. Also, clear comparisons and benchmarks between RL and other design methods are lacking in previous literature. Therefore, at the moment, the suitability of RL for process design is questionable and remains an open question. Here, we present the main differences between RL and existing superstructure-free methods in the context of process design:

- **Computational efficiency (static vs. dynamic optimization)** - RL is better suited for dynamic, sequential-decision problems rather than static ones. However, prior studies have applied RL to stationary process design, constituting static optimization problems. This may drastically increase complexity and decrease computational efficiency (called "sample efficiency" in RL). As a side note, the capability of RL in dynamic optimization paves the way for solving integrated design and operation problems (see Section 3.5).
- **Iterative build-up vs. full flowsheet generation** - RL sequentially generates flowsheets (unit by unit), differing from evolutionary superstructure-free methods that construct flowsheets as a whole. While this sequential build-up increases computational complexity, there are also potential advantages. Generating feasible flowsheets and simulating them is challenging, often causing convergence problems in evolutionary superstructure-free methods. In contrast, the iterative strategy of RL promotes convergence and intermediate simulations of incomplete flowsheets may provide valuable information for learning. However, this additional information comes at the cost of additional simulation time.
- **Inference (online) vs. optimization (offline)** - A significant distinction between RL and other optimization methods lies in the solution time for similar, recurring problems. RL involves training a policy once, which is then utilized during inference to rapidly predict near-optimal solutions, offering a significant advantage in time-sensitive control applications. In contrast, classical optimization algorithms solve problems individually, typically requiring long runtimes for each problem instance. In process design, long optimization times are usually not an issue (unless it becomes intractable

in the case of large non-convex MINLPs). Thus, classical optimization methods are usually well-suited. However, the rapid inference capability of RL may also provide new opportunities for process design. For example, RL agents might be integrated into flowsheet simulation software to automatically and immediately suggest near-optimal design options to users. Also, fast solutions of design problems can be advantageous when a large numbers of design problems need to be solved (e.g., as subproblems in larger optimization studies).

- **Learning capacity** - RL possesses a substantially greater learning capacity compared to standard evolutionary methods (e.g., more trainable parameters). For instance, state-of-the-art deep RL algorithms can incorporate extensive networks of learnable parameters, potentially exceeding billions of parameters. This enhanced learning potential facilitates inference and allows for retaining information, unlike genetic algorithms which typically lose details about previous populations. This substantial learning capacity of RL holds the potential for learning more complex dependencies between design actions and results. However, the high learning capacity also presents severe drawbacks such as a vast amount of training data and extensive training duration (measured in epochs within RL).

When comparing RL with standard optimization methods, its most notable advantages include larger learning capacity and inference ability. However, at the same time, RL typically demands significantly more training simulations than static optimization solvers, which leaves doubt on whether its potential advantages outweigh the disadvantages in the context of process design. Current literature applying RL to process design neglects its inference capabilities, learning capacity, and dynamic optimization capabilities, predominantly utilizing the training phase of RL as an evolutionary optimization strategy for static problems. Thereby, they essentially merge the drawbacks of both worlds. Additionally, there is a lack of computational comparison between RL and traditional process design methods. In the following, we critically examine the existing literature on RL in process design (Section 2) and highlight future perspectives (Section 3).

2 State of the art

The general framework of RL in process design is shown in Figure 2. The agent learns to design processes by iteratively placing unit operations with design and operating variables, and simulating the resulting processes in the environment, ultimately obtaining the optimal policy π^* which designs optimal processes. Mathematically, this problem can be formulated as Markov decision processes (MDP): $M = \{S, A, T, R\}$ with states $\mathbf{s} \in S$, actions $\mathbf{a} \in A$, the transition function $T : S \times A \times S \rightarrow [0, 1]$, and the reward function $R : S \times A \times S \rightarrow \mathbb{R}$. In the context of process design, the states \mathbf{s} represent the flowsheet topology as well as all relevant design specifications, operating variables, thermodynamic stream data, flowrates,

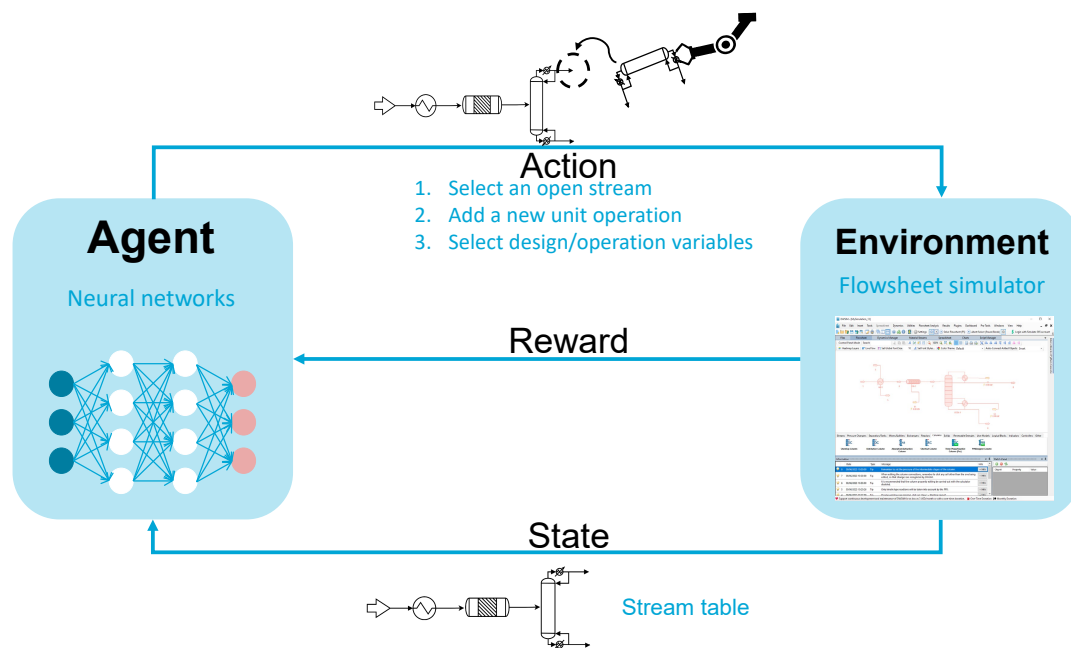


Figure 2: General framework of reinforcement learning for process design.

and compositions. The agent takes the current states \mathbf{s} as input to take actions \mathbf{a} . These actions can include design and operating variables. In chemical processes, design variables are usually determined during the initial design phase and typically remain fixed throughout the operation, such as equipment size. Operating variables can be adjusted during operation (e.g., flow rates, and pressures). Furthermore, the actions contain discrete choices (e.g., the selection of open streams, unit operation types, or number of stages) as well as continuous choices (e.g., the length of a reactor or operating flowrates), namely hybrid action space. Usually, decisions in process design are also hierarchical. For example, the agent first determines an open stream to add a unit operation, then the type of unit operation, then design variables, and finally operating variables. After a new unit operation is added, the new flowsheet is simulated in an environment (e.g., a process simulation software). After finishing a flowsheet, a numerical reward \mathbb{R} is returned to the agent. This corresponds to the objective for optimization. By repeating the design of multiple flowsheets and receiving corresponding rewards, the agent learns to design processes that maximize the reward.

In this section, we survey the RL state-of-the-art literature (summarized in Table 1) based on (i) information representation, (ii) agent architecture, and (iii) environment and reward.

2.1 Information representation

Chemical processes comprise various information, such as process topology, thermodynamic states, flowrates, concentrations, design variables, operating variables, components, and underlying mechanistic knowledge. The meaningful representation of the chemical process information is critical for the learning

Table 1: An overview of the reviewed literature and different choices of elements in RL for process design. We use Repr. to indicate information representation. Furthermore, we utilize Dis., Cont., and Hier. to denote discrete, continuous, and hierarchical action space of RL. Within the decisions of RL, we use Topo., Des. and Oper. to represent actions involved in changing flowsheet topologies, selecting design variables and operating variables, respectively.

Ref.	Repr.	Agent architecture	Environment	Action space			Decisions		
				Dis.	Cont.	Hier.	Topo.	Des.	Oper.
[22]	Matrix	Actor-Critic (SAC)	COCO	✓	✓		✓	✓	✓
[23]	Matrix	Actor-Critic (SAC)	Aspen Plus	✓	✓			✓	✓
[18]	Matrix	Actor-Critic (-)	Aspen Plus		✓				✓
[30]	Matrix	Actor-Critic (PPO)	Short-cut		✓			✓	
[24] [25] [26]	Matrix	Actor-Critic (Two players)	Short-cut	✓		✓	✓		
[28][29]	Matrix	Actor-Critic (PPO)	Short-cut	✓	✓	✓	✓	✓	✓
[20]	Matrix	Actor-Critic (ACER)	Short-cut	✓				✓	
[19]	Matrix	Policy-based (Policy search)	Short-cut		✓				✓
[21]	Matrix	Policy-based (PG)	Short-cut		✓			✓	✓
[32]	Matrix	Value-based (DQN)	IDAES	✓			✓		
[31]	Matrix	Value-based (DQN)	UniSim		✓			✓	✓
[27]	Matrix	Value-based (Q-learning)	Short-cut	✓	✓		✓	✓	
[33]	Graph	Actor-Critic (PPO)	Short-cut	✓	✓	✓	✓	✓	✓
[34]	Graph	Actor-Critic (PPO)	DWSIM	✓	✓	✓	✓	✓	✓

and generalization of RL agents. For RL in process design, there are currently two methods for information representations: Matrix [22, 19, 23, 30, 24, 25, 26, 20, 29, 21, 31, 18, 32, 27, 28] and graph [33, 34].

In matrix-based representation, flowsheets are represented by fixed-size matrices. Within the flowsheet matrix, the connectivity, stream compositions, thermodynamic stream data, and design variables are usually concatenated. For example, Göttl et al. [24, 25, 26] represented flowsheets as 16×28 matrices, where each row represents a stream and encompasses four parts: $\{\mathbf{v}, \mathbf{u}, \mathbf{d}, \mathbf{t}\}$. \mathbf{v}_i has five entries, which describe the molar fractions and total molar flowrate of stream i . \mathbf{u}_i stores the type of the unit operation that is downstream of stream i as one-hot encoding. Furthermore, \mathbf{d}_i stores the connectivity of unit operations and has sixteen entries (i.e., this corresponds to the adjacency matrix). Finally, \mathbf{t}_i has two entries: the first entry indicates whether the task is terminated (0 if not terminated), and the second entry indicates whether stream i is still unused (0 if unused). Most of the previous publications use matrix representations of flowsheet states (c.f. Table 1).

In graph-based representation, flowsheets are represented by directed heterogeneous graphs. Flowsheet graphs consist of nodes and edges. Unit operations are represented by nodes, also referred to as vertices $v \in V$, and streams are represented by edges $e_{vw} \in E$ connecting two nodes v and w . Importantly, node feature vector $\mathbf{f}^v \in F^V$, and edge feature vector $\mathbf{f}^{e_{vw}} \in F^E$ are associated with each node and edge, respectively. Within the node feature vectors, types of unit operation, design specifications, and operating points are encoded. The edge feature vectors contain thermodynamic states, concentrations, and flowrates. In the past, only our previous works used graph representations of flowsheets for RL [33, 34].

The comparison between flowsheet matrices and flowsheet graphs is still an open research question in

the context of RL in process design. Flowsheet matrices are easier to implement than flowsheet graphs and are used by the majority of the literature as shown in Table 1. Flowsheet matrices are processed by RL agents using multilayer perceptrons (MLPs) or convolutional neural networks (CNNs). However, every flowsheet graph has $N!$ different adjacency matrices. CNNs and MLPs are not permutation equivariance

$$f(\mathbf{P}^T \mathbf{x} \mathbf{P}) \neq \mathbf{P}^T f(\mathbf{x}) \mathbf{P}$$

where \mathbf{P} is a permutation matrix, \mathbf{x} is the input matrix and f is a MLP or CNN [35]. This means that such models depend on the arbitrary order of rows/columns in the flowsheet matrix and thus, cannot generalize over flowsheet topologies. Also, the neighborhood of an entry in the matrix does not correspond to physical connectivity which makes learning using MLPs/CNNs more difficult as it requires learning long-range interactions. In contrast, (message-passing) graph convolutional networks (GCNs) are permutation equivariance and they learn from the actual connectivity of flowsheet graphs. Furthermore, MLPs and CNNs require fixed-size inputs, e.g., a pre-defined maximum number of unit operations and streams, while GCNs are size-independent [36].

2.2 Agent architecture

RL agents consist of two components: A policy and a learning algorithm. The policy describes the behavior of the agent, mapping the current state \mathbf{s} into an action \mathbf{a} : $\pi(\mathbf{s}) = \mathbf{a}$. It is parameterized by function approximators such as MLPs. The learning algorithm is used to continuously update the policy based on the actions, states, and rewards. Depending on the learning algorithms, the agent can be characterized into three types: Value-based, policy-based, and actor-critic-based (AC).

Value-based agents learn a functional approximator of the value function ($V_\pi(\mathbf{s})$) to take actions. The value function outputs the expected returns after the current process step t given a state \mathbf{s} and a policy π : $V_\pi(\mathbf{s}) = \mathbb{E}_\pi [G_t | \mathbf{s}]$, where returns G_t :

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^k R_{t+k+1} = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (1)$$

where γ^k are the discount rates to determine the present value of future rewards, and k is the process step from t to the end of the episode. Similarly, we can also derive the state-action value function, namely the quality function Q_π , which calculates the expected returns given a state \mathbf{s} and action \mathbf{a} , following policy π : $Q_\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E}_\pi [G_t | \mathbf{s}, \mathbf{a}]$. Depending on the calculated V-/Q-value, different search algorithms such as best-first search or nearest neighbors are used to choose the final action. In the context of RL in process design, three works [27, 32, 31] deployed Q-learning based agent to perform process synthesis tasks. However, traditional value-based agents can only take discrete actions, which

hinders further development because continuous decision-making of design or operating variables is vital in process design tasks.

Policy-based agents directly learn a functional approximator of the policy function. Specifically, the policy approximator π_θ maps the current states \mathbf{s} to the actions \mathbf{a} : $\pi_\theta(\mathbf{s}) = \mathbf{a}$. And the optimal policy π^* is obtained by maximizing the expected return $\mathbb{E}_\theta [G_t]$ through policy gradient or policy-search methods. In the context of RL in process design, Sachio et al. [21] and Perera et al. [19] utilized policy gradient methods and policy-search methods to perform process design tasks, respectively. Compared to the value-based approach, the policy-based agent can handle both discrete and continuous actions. However, policy-based methods are known for high variance and sub-optimal local solutions [37].

AC agents combine the advantages of value-based and policy-based methods. AC consists of an actor, working as a functional approximator of the policy function, and a critic, serving as a functional approximator of the value function. Therefore, AC agents explicitly optimize both value and policy functions and are able to process both discrete and continuous action spaces. In the context of RL in process design, different types of AC agents have been used such as Proximal Policy Optimization (PPO) [29] [30] [33] [34], Soft Actor-Critic (SAC) [22, 23], Two-player game [24, 25, 26], and Sample Efficient Actor Critic with Prioritized Experience Replay (ACER) [20].

The choice of agent architecture for RL in process design is an open question. AC RL is deemed to be a viable option because it combines the advantages of value-based and policy-based and can handle both discrete and continuous decisions. Specifically, PPO is the most popular algorithm in process design tasks with the advantage of less complicated implementation and a stable learning process. However, PPO is an on-policy algorithm which means the optimized policy is the same as the policy for action selection. Therefore, PPO is less data-efficient than off-policy algorithms, such as SAC and ACER, which may take less time and fewer training episodes. Moreover, AC RL comes with challenges, including complex implementation, computational demands when optimizing both actor and critic networks concurrently, and potential convergence issues [37, 13].

2.3 Environment and reward

The environment simulates the processes and computes a reward as feedback to the agent. Selecting an appropriate accuracy level for the environment is a vital task for RL in process design and depends on the task-specific requirements and modeling intent. There are two main levels of accuracy: Shortcut and rigorous simulators. Shortcut simulators utilize approximated process models to ensure tractability but can be inaccurate. Rigorous simulators involve more accurate process models that require longer computation times, as previous studies indicated [34, 23]. In the past, RL for process design has used multiple process simulation software including open-source (DWSIM, IDEAS), non-commercial (COCO), and commercial (UniSim, Aspen Plus) alternatives. Additionally, Seidenberg et al. [29] leveraged knowl-

edge graphs to retrieve information about the design task, process knowledge, and the current state of the process. Notably, this knowledge graph was part of a manually-implemented environment and not directly accessible to the RL agent. Thus, the agent also relied on a flowsheet matrix representation as states.

Previous research optimized towards a single economic objective [22, 23, 29, 24, 25, 26, 20, 33, 34, 27, 21, 18]. Also, some works integrate purity, recovery, power consumption, and product flow rate into scalar reward functions [30, 32, 31].

3 Perspectives

Despite the first demonstrations of RL for process design, it is still unclear if RL outperforms existing design methods. In our view, the big research challenge is the generalization of RL models and the use of its inference capabilities. The training phase of current RL frameworks is essentially used like a derivative-free optimization approach (e.g., a genetic algorithm) to optimize the process topology for one particular case study. Thus, a re-training is needed for a new case study and the agent fails to transfer its learning to new situations. In general, deep RL has an inference and a high learning capacity. The derivative-free optimization approach with RL does not use the full potential of RL. However, useful application of inference requires generalization across multiple case studies. This generalization requires an extension of the information representation and agent architecture to account for process-relevant knowledge. This includes domain expertise, prior process data, and physical constraints which are typically employed by engineers when designing chemical processes. Integrating this information would allow the RL agent to see what "drives the process" and ultimately unlock the full potential of RL by learning from multiple processes. We envision that RL will generalize (to some extent) and ultimately design processes at inference time. In this section, we provide our perspectives on underlying challenges and a number of other promising future works.

3.1 Information representation

Information representation is critical for RL since it encapsulates the current state of the environment, which directly affects decision-making for agents. However, current information representations still lack mechanistic knowledge and relevant process information. Furthermore, it neglects the appropriate representation of molecules. Integrating the above information will significantly benefit the RL agent in generalization. Numerous representation methods could be potentially incorporated into RL in process design for process and molecular information representation. For example, Simplified Flowsheet Input-Line Entry-System (SFILES) [38], SFILES 2.0 [39], eSFILES [40], and knowledge graphs [29] could be used to enhance process representations. Similarly, molecular descriptors [41], molecular graphs [42],

SMILES [43], knowledge graphs [44], and hypergraphs [45] could be used to encode molecular information.

3.2 Agent architecture

In this section, we identify the limitations and potential improvements of the current agent architecture.

3.2.1 Integration of mechanistic knowledge

Current RL algorithms are not sufficient to transfer knowledge into the new processes because RL agents have a limited understanding of mechanistic knowledge and physical properties. Future work could consider implementing a physics-informed RL agent by encoding information-rich representations such as knowledge graphs or hypergraphs to inform the agent. Furthermore, fundamental concepts, such as thermodynamic driving forces (Gibb’s free energy), could be included in the RL agents. This allows the agent to learn general concepts that can be translated into other problems because they are based on physics.

3.2.2 Integration of prior data

RL for process design is currently initialized randomly, which can lead to suboptimal solutions, excessive training times, and frequent convergence issues. Meanwhile, there is a large number of existing digitized chemical process data from simulation files and images [46], which can potentially accelerate the learning process of the agent. Transfer learning improves learning performance by transferring knowledge from different but relevant domains [47]. In RL for process design, three work [34, 32, 21] already leveraged transfer learning to accelerate the learning process, e.g., from short-cut simulators to rigorous simulators [34] and from one case study to another case study [32]. However, in the current transfer learning setting, the agent is still not learning from existing chemical process information. Future work can consider leveraging encoder-decoder models such as Variational Autoencoders (VAEs) or transformers to learn from existing flowsheets and then applying transfer learning to the agent.

3.2.3 Stochastic decision-making

Considering the uncertainty of energy/feedstock prices and demand is a major challenge for renewable processes [7]. However, current RL agents ignore fluctuations in energy/feedstock prices and demand. Future research could separate design and operating variables in the RL agent. This allows the inclusion of multiple scenarios for flexible operation. Besides, additional encoders or actors can be included to process stochastic energy prices, demands, and raw material compositions as additional inputs at an operational level. Therefore, the agent can automatically select the operating variables based on stochastic energy prices and demand in two-stage stochastic programming settings.

3.2.4 Constrained decision-making

Constrained decision-making is crucial for RL in process design to ensure optimal and safe performance. However, standard RL agent frameworks cannot enforce constraints but include constraints as "soft" penalties in the reward functions [33, 34, 23]. Future work should focus on integrating constraints directly in the agent structure. As a first step, an additional critic network could be built to account for safety constraints, guiding RL agents to explore appropriate regions in policy optimization [48].

3.3 Environment and reward

In this section, we offer our perspectives on the limitations of the environment and reward setup and provide several suggestions for future work.

3.3.1 Standardized simulation interfaces

RL agents frequently interact with process simulators during the training process and the interaction relies on individual interfaces as Table 1 shows. However, current interfaces are usually simulator-specific, which means that a new interface needs to be implemented from scratch whenever a new process simulator is included. This process is highly repetitive and inefficient, especially for incorporating multi-fidelity process models. Future work could implement a standardized simulation interface that enables the agent to exchange data efficiently and uniformly between different process simulators. This interface could potentially make use of existing standards such as CAPE-OPEN [49] and DEXPI+ [50].

3.3.2 Multi-fidelity process models

Current research only leverages a single fidelity model for RL in process design tasks. However, RL agents greatly benefit from pre-training on low-fidelity process simulators and subsequent fine-tuning on high-fidelity process simulators [34]. Therefore, future research can focus on developing an agent that can dynamically select between multiple fidelity models during training. Specifically, a probabilistic model can be developed to guide the RL actor based on multi-fidelity critics to reduce training times and resolve convergence issues.

3.3.3 Multi-objective rewards

Current RL frameworks for process design are not suitable for sustainable process design because they are limited to a single objective function. In the future, the current agent structures could be extended to include multiple objectives. For example, the critic network could predict multiple rewards, which will be processed by multi-objective optimization to generate corresponding weights for each objective (e.g. economic, environmental, safety) [51].

3.4 Integrated molecular and process design

Current RL frameworks lack the co-design of molecules. However, the design or selection of molecules is a critical task in many process design tasks, e.g., co-design of working fluids, solvents, or products [52]. Also, RL has already been used for molecular design [53]. Therefore, future work should consider integrating these concepts using RL. For instance, future work could first use an RL agent for molecular design (e.g., based on [53]) to design a solvent. Then, a mechanistic or data-driven model [42] can be used to estimate the relevant properties of the generated molecules. Subsequently, the properties are utilized to simulate the process within the RL for the process design framework. This essentially adds a new hierarchy level to the existing RL for the process design framework.

3.5 Integrated process operation and design

Integrating process design and process control becomes increasingly relevant as renewable energy and feedstock demand fluctuates. For example, the design of a process could be optimized while simultaneously optimizing its operation under changing feedstock compositions or energy prices. Another example is the optimal design of batch processes while considering optimal operation strategies. These problems are usually formulated as mixed-integer dynamic optimization (MIDO) problems which are difficult to solve. RL is a tool to solve discrete-time dynamic optimization problems, which makes it suitable for integrating process operation and design. However, current RL works only consider process design and operation separately or focus on a specific unit operation design and operation [21]. Future research could integrate process design with process operation through the RL. For instance, future work could extend the hierarchical RL agent into separate design and operation agents. Then, the design agent defines the design variables, and the operation agent subsequently optimizes operating variables given the current design. Notably, this would require the use of a dynamic simulation environment and will lead to high computational demands. Thus, future research is needed to solve the resulting multi-scale problem efficiently.

3.6 Benchmarking with established methods

Numerous established methods are available to solve design optimization problems, including deterministic (e.g., BARON or MAiNGO) and stochastic solvers (e.g., genetic algorithms, Bayesian optimization) [54]. It is still questionable how RL compares against these traditional approaches for steady-state process design. Dynamic solution approaches such as RL can in principle be used to solve static optimization problems but are likely significantly less efficient. However, many process design problems in chemical engineering are actually (mixed-integer) dynamic optimization problems (c.f. Section 3.5). In such instances, RL may be an efficient solution approach.

Future work should carefully assess the advantages and disadvantages of using RL for steady-state process design and static optimization in general. We recommend conducting comparisons to benchmark different methods. Moreover, we envision the development of new ML-based algorithms that integrate some of the advantages of RL (e.g., large learning capacity and inference capabilities) in the context of process design. For example, encoder-decoder models could be combined with active learning to predict process flowsheet graphs directly [8].

4 Conclusions

We reviewed the state-of-the-art RL in process design in terms of information representation, agent architecture, environment, and reward. RL has shown initial promising results for process design but its suitability in static process design still needs to be examined. Additionally, a detailed comparison with existing process design methods is missing and current RL frameworks show limited generalization capabilities. Therefore, we advocate that future research should benchmark RL with other process design methods. Additionally, to unlock the full potential of RL, new concepts for meaningful information representation are required. Furthermore, the integration of mechanistic knowledge, existing process data, uncertainties, and constraints would be highly beneficial for optimal decision-making. Finally, future RL frameworks could also integrate molecular design and process operation into the conceptual process design.

References

- [1] Samir Isaac Meramo-Hurtado and Ángel Dario González-Delgado. Process synthesis, analysis, and optimization methodologies toward chemical process sustainability. *Industrial & Engineering Chemistry Research*, 60(11):4193–4217, 2021.
- [2] Elias Martinez-Hernandez. Trends in sustainable process design—from molecular to global scales. *Current Opinion in Chemical Engineering*, 17:35–41, 2017.
- [3] Tomio Umeda. Computer aided process synthesis. *Computers & Chemical Engineering*, 7(4):279–309, 1983.
- [4] Alexandre C Dimian and Costin Sorin Bildea. *Chemical process design: Computer-aided case studies*. John Wiley & Sons, 2008.
- [5] T.F. Yee and I.E. Grossmann. Simultaneous optimization models for heat integration—II. heat exchanger network synthesis. *Computers & Chemical Engineering*, 14(10):1165–1184, 1990.

- [6] Luca Mencarelli, Qi Chen, Alexandre Pagot, and Ignacio E. Grossmann. A review on superstructure optimization approaches in process system engineering. *Computers & Chemical Engineering*, 136:106808, 2020.
- [7] Alexander Mitsos, Norbert Asprion, Christodoulos A. Floudas, Michael Bortz, Michael Baldea, Dominique Bonvin, Adrian Caspari, and Pascal Schäfer. Challenges in process optimization for new feedstocks and energy sources. *Computers & Chemical Engineering*, 113:209–221, 2018.
- [8] Edwin Hirtreiter, Lukas Schulze Balhorn, and Artur M. Schweidtmann. Toward automatic generation of control structures for process flow diagrams with large language models. *AIChE Journal*, 2023.
- [9] Tahar Nabil, Jean-Marc Commenge, and Thibaut Neveux. Generative approaches for the synthesis of process structures. In *Computer Aided Chemical Engineering*, volume 49, pages 289–294. Elsevier, 2022.
- [10] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning, 2013.
- [11] Michal Kempka, Marek Wydmuch, Grzegorz Runc, Jakub Toczek, and Wojciech Jaskowski. ViZDoom: A doom-based AI research platform for visual reinforcement learning. In *2016 IEEE Conference on Computational Intelligence and Games (CIG)*. IEEE, 2016.
- [12] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, Timothy Lillicrap, Karen Simonyan, and Demis Hassabis. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144, 2018.
- [13] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [14] J.C. Hoskins and D.M. Himmelblau. Process control via artificial neural networks and reinforcement learning. *Computers & Chemical Engineering*, 16(4):241–251, 1992.
- [15] S.P.K. Spielberg, R.B. Gopaluni, and P.D. Loewen. Deep reinforcement learning approaches for process control. In *2017 6th International Symposium on Advanced Control of Industrial Processes (AdCONIP)*. IEEE, 2017.
- [16] Christian D. Hubbs, Can Li, Nikolaos V. Sahinidis, Ignacio E. Grossmann, and John M. Wassick. A deep reinforcement learning approach for chemical production scheduling. *Computers & Chemical Engineering*, 141:106982, 2020.

- [17] Young Hoon Lee and Seunghoon Lee. Deep reinforcement learning based scheduling within production plan in semiconductor fabrication. *Expert Systems with Applications*, 191:116222, 2022.
- [18] Junhui Chen and Fan Wang. Cost reduction of CO₂ capture processes using reinforcement learning based iterative design: A pilot-scale absorption–stripping system. *Separation and Purification Technology*, 122:149–158, 2014.
- [19] A.T.D. Perera, P.U. Wickramasinghe, Vahid M. Nik, and Jean-Louis Scartezzini. Introducing reinforcement learning to the energy system design process. *Applied Energy*, 262:114580, 2020.
- [20] Cesare Caputo, Michel-Alexandre Cardin, Pudong Ge, Fei Teng, Anna Korre, and Ehecatl Antonio del Rio Chanona. Design and planning of flexible mobile micro-grids using deep reinforcement learning. *Applied Energy*, 335:120707, 2023.
- [21] Steven Sachio, Max Mowbray, Maria M. Papathanasiou, Ehecatl Antonio del Rio-Chanona, and Panagiotis Petsagkourakis. Integrating process design and control using reinforcement learning. *Chemical Engineering Research and Design*, 183:160–169, 2022.
- [22] Laurence Illing Midgley. Deep reinforcement learning for process synthesis. 2020.
- [23] Stephan C. P. A. van Kalmthout, Laurence I. Midgley, and Meik B. Franke. Synthesis of separation processes with reinforcement learning. 2022.
- [24] Quirin Göttl, Dominik G. Grimm, and Jakob Burger. Automated synthesis of steady-state continuous processes using reinforcement learning. *Frontiers of Chemical Science and Engineering*, 16(2):288–302, 2021.
- [25] Quirin Göttl, Yannic Tönges, Dominik G. Grimm, and Jakob Burger. Automated flowsheet synthesis using hierarchical reinforcement learning: Proof of concept. *Chemie Ingenieur Technik*, 93(12):2010–2018, 2021.
- [26] Quirin Göttl, Dominik G. Grimm, and Jakob Burger. Using reinforcement learning in a game-like setup for automated process synthesis without prior process knowledge. In *Computer Aided Chemical Engineering*, pages 1555–1560. Elsevier, 2022.
- [27] Ahmad Khan and Alexei Lapkin. Searching for optimal process routes: A reinforcement learning approach. *Computers & Chemical Engineering*, 141:107027, 2020.
- [28] Ahmad A. Khan and Alexei A. Lapkin. Designing the process designer: Hierarchical reinforcement learning for optimisation-based process design. *Chemical Engineering and Processing - Process Intensification*, 180:108885, 2022.

- [29] J. Raphael Seidenberg, Ahmad A. Khan, and Alexei A. Lapkin. Boosting autonomous process design and intensification with formalized domain knowledge. *Computers & Chemical Engineering*, 169:108097, 2023.
- [30] Siby Jose Plathottam, Blake Richey, Gregory Curry, Joe Cresko, and Chukwunwike O. Iloeje. Solvent extraction process design using deep reinforcement learning. *Journal of Advanced Manufacturing and Processing*, 3(2), 2021.
- [31] Sam Kim, Mun-Gi Jang, and Jin-Kuk Kim. Process design and optimization of single mixed-refrigerant processes with the application of deep reinforcement learning. *Applied Thermal Engineering*, 223:120038, 2023.
- [32] Dewei Wang, Jie Bao, Miguel Zamarripa-Perez, Brandon Paul, Yunxiang Chen, Peiyuan Gao, Tong Ma, Alexander Noring, Arun Iyengar, Daniel Schwartz, Erica Eggleton, Qizhi He, Andrew Liu, Olga Marina, Brian Koeppe, and Zhijie Xu. Reinforcement learning for automated conceptual design of advanced energy and chemical systems. 2022.
- [33] Laura Stops, Roel Leenhouts, Qinghe Gao, and Artur M. Schweidtmann. Flowsheet generation through hierarchical reinforcement learning and graph neural networks. *AIChE Journal*, 69(1), 2022.
- [34] Qinghe Gao, Haoyu Yang, Shachi M. Shanbhag, and Artur M. Schweidtmann. Transfer learning for process design with reinforcement learning. In *Computer Aided Chemical Engineering*, pages 2005–2010. Elsevier, 2023.
- [35] William L. Hamilton. *Graph Representation Learning*. Springer International Publishing, 2020.
- [36] Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. Graph neural networks: A review of methods and applications. *AI Open*, 1:57–81, 2020.
- [37] Ofir Nachum, Mohammad Norouzi, Kelvin Xu, and Dale Schuurmans. Bridging the gap between value and policy based reinforcement learning. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [38] Loïc d’Anterrockes. *Process Flowsheet Generation & Design through a Group Contribution Approach*. [CAPEC], Department of Chemical Engineering, Technical University of Denmark, 2005.
- [39] Gabriel Vogel, Edwin Hirtreiter, Lukas Schulze Balhorn, and Artur M. Schweidtmann. SFILES 2.0: an extended text-based flowsheet representation. *Optimization and Engineering*, 2023.

- [40] Vipul Mann, Rafiqul Gani, and Venkat Venkatasubramanian. Intelligent process flowsheet synthesis and design using extended SFILES representation. In *Computer Aided Chemical Engineering*, pages 221–226. Elsevier, 2023.
- [41] Roberto Todeschini and Viviana Consonni. Molecular descriptors. *Recent Advances in QSAR Studies*, pages 29–102, 2010.
- [42] Artur M. Schweidtmann, Jan G. Rittig, Andrea König, Martin Grohe, Alexander Mitsos, and Manuel Dahmen. Graph neural networks for prediction of fuel ignition quality. *Energy & Fuels*, 34(9):11395–11407, 2020.
- [43] David Weininger. SMILES, a chemical language and information system. 1. introduction to methodology and encoding rules. *Journal of Chemical Information and Computer Sciences*, 28(1):31–36, 1988.
- [44] Yin Fang, Qiang Zhang, Haihong Yang, Xiang Zhuang, Shumin Deng, Wen Zhang, Ming Qin, Zhuo Chen, Xiaohui Fan, and Huajun Chen. Molecular contrastive learning with chemical element knowledge graph. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(4):3968–3976, 2022.
- [45] Hiroshi Kajino. Molecular hypergraph grammar with its application to molecular optimization. 2018.
- [46] Lukas Schulze Balhorn, Qinghe Gao, Dominik Goldstein, and Artur M. Schweidtmann. Flowsheet recognition using deep convolutional neural networks. In *Computer Aided Chemical Engineering*, pages 1567–1572. Elsevier, 2022.
- [47] Karl Weiss, Taghi M. Khoshgoftaar, and DingDing Wang. A survey of transfer learning. *Journal of Big Data*, 3(1), 2016.
- [48] Qisong Yang, Thiago D. Simão, Simon H. Tindemans, and Matthijs T. J. Spaan. Safety-constrained reinforcement learning with a distributional safety critic. *Machine Learning*, 112(3):859–887, 2022.
- [49] M Jarke, J Koeller, W Marquardt, L von Wedel, and B Braunschweig. Cape-open: Experiences from a standardization effort in chemical industries. In *Proc. of 1st IEEE Conference on Standardisation and Innovation in Information Technology (SIIT 99)(Aachen, Germany)*, pages 25–35, 1999.
- [50] M Theissen and M Wiedau. Dexpi p&id specification. *Version 0.11*, 2016.
- [51] Chunming Liu, Xin Xu, and Dewen Hu. Multiobjective reinforcement learning: A comprehensive overview. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 45(3):385–398, 2015.

-
- [52] Philipp Rehner, Johannes Schilling, and André Bardow. Molecule superstructures for computer-aided molecular and process design. *Molecular Systems Design & Engineering*, 8(4):488–499, 2023.
- [53] Marcus Olivecrona, Thomas Blaschke, Ola Engkvist, and Hongming Chen. Molecular de-novo design through deep reinforcement learning. *Journal of Cheminformatics*, 9(1), 2017.
- [54] Fani Boukouvala, Ruth Misener, and Christodoulos A. Floudas. Global optimization advances in mixed-integer nonlinear programming, MINLP, and constrained derivative-free optimization, CDFO. *European Journal of Operational Research*, 252(3):701–727, 2016.