



Faculty of Computer Science

Study Program Placeholder

PRISMA-Guided Reinforcement Learning for Manufacturing Scheduling

Master Thesis

von

OpenAI Codex

Submission Date: dd.mm.yyyy

First Examiner: Prof. Dr. Noah Klarmann

Second Examiner: Prof. Dr. Placeholder

EIGENSTÄNDIGKEITSERKLÄRUNG / DECLARATION OF ORIGINALITY

Hiermit bestätige ich, dass ich die vorliegende Arbeit selbständig verfasst und keine anderen als die angegebenen Hilfsmittel benutzt habe. Die Stellen der Arbeit, die dem Wortlaut oder dem Sinn nach anderen Werken (dazu zählen auch Internetquellen) entnommen sind, wurden unter Angabe der Quelle kenntlich gemacht.

I declare that I have authored this thesis independently, that I have not used other than the declared sources / resources, and that I have explicitly marked all material which has been quoted either literally or by content from the used sources.

Rosenheim, den November 10, 2025

OpenAI Codex

Abstract

This thesis presents a PRISMA-compliant literature review on reinforcement learning for manufacturing scheduling, supported by a multi-agent LLM workflow. Replace this placeholder with the finalized abstract.

Keywords: reinforcement learning, manufacturing scheduling, PRISMA, multi-agent systems, literature review

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Research Questions	1
1.3	Scope and Assumptions	1
1.4	Contributions	2
1.5	Industrial Scope and Constraints	2
1.6	Industry Interviews and Practitioner Pain Points	2
1.7	Terminology and Definitions	2
1.8	Thesis Structure	3
2	Background	4
2.1	Manufacturing Scheduling Fundamentals	4
2.2	Reinforcement Learning Primer	4
2.3	Evaluation Metrics	5
2.4	Digital Twins and Data Infrastructure	5
2.5	Benchmark Instances and Datasets	6
2.6	Operations Research Synergies	6
2.7	Human Factors and Safety Considerations	6
3	Methodology	8
3.1	PRISMA Workflow	8
3.1.1	Information Sources	8
3.1.2	Eligibility Criteria	8
3.1.3	PRISMA Reporting	8
3.2	Data Extraction	8
3.3	Multi-Agent Process	9
3.4	Manufacturing-Focused Screening Details	10
3.5	Data Quality Assurance and Tooling	10
3.6	Search String Construction and Validation	11
3.7	Risk of Bias and Quality Assessment	11
3.8	Qualitative Coding of Findings	11
3.9	Data Refresh Cadence	11
4	Experimental Protocols and Evaluation Practices	12
4.1	Simulation Setup Patterns	12
4.2	Training Pipelines	12
4.3	Baseline Selection	13
4.4	Statistical Validation and Uncertainty	13
4.5	Benchmarking Limitations	13

5	Toolchain and Automation	14
5.1	Repository Layout	14
5.2	Automation Scripts	14
5.3	Line and Bar Plot Generation	15
5.4	Validation and Logging	15
5.5	Reproducible Builds	15
5.6	Extensibility	15
5.7	Collaboration Workflow	16
5.8	Limitations and Future Tooling	16
6	Literature Landscape	17
6.1	Descriptive Statistics	17
6.2	RL Technique Distribution	18
6.3	Domain Spotlights	18
6.3.1	Job-Shop and Flexible Cells	18
6.3.2	Flow-Shop, Hybrid, and Battery Lines	19
6.3.3	Semiconductor Ecosystem	19
6.3.4	Energy, Microgrid, and Sustainability	19
6.3.5	Specialized Manufacturing Domains	20
6.4	Technology Readiness Indicators	20
6.5	Data Availability Patterns	21
6.6	Scheduling Objectives	21
7	Industrial Case Studies	26
7.1	Flexible Job Shops and Aerospace Cells	26
7.2	Flow Shops, Battery Lines, and Hybrid Plants	26
7.3	Semiconductor Manufacturing and Supply Chains	27
7.4	Energy-Aware and Microgrid-Integrated Factories	27
7.5	Regulated Industries: Pharma, Biopharma, and Remanufacturing	28
7.6	Circular Manufacturing and Long-Horizon Planning	28
7.7	Cross-Domain Lessons	28
8	Comparative Analysis	30
8.1	Job-Shop and Flexible Manufacturing Cells	30
8.2	Flow-Shop, Hybrid, and Aerospace Lines	30
8.3	Semiconductor and High-Mix Electronics	31
8.4	Energy, Microgrid, and Sustainability-Oriented Scheduling	31
8.5	Specialized Cells: Robotics, Pharma, and Human-Centric Lines	32
8.6	Interpretability and Human-in-the-Loop Governance	32
9	Quantitative Meta-Analysis	33
9.1	Makespan and Tardiness Improvements	33
9.2	Energy and Carbon Reductions	33
9.3	Statistical Rigor	33
9.4	Sensitivity to Reward Design	34
9.5	Deployment Readiness Scorecard	34
9.6	Data Availability Metrics	34

9.7	Implications	35
10	Roadmap to a 200-Study Corpus	36
10.1	Target Domains and Venues	36
10.2	Search Automation Enhancements	36
10.3	Inclusion of Non-English Sources	36
10.4	Data Harmonization	36
10.5	Timeline and Milestones	37
10.6	Community Engagement	37
11	Slidex Communication Layer	38
11.1	Design Principles	38
11.2	Data Synchronization	38
11.3	Speaker Notes and Q&A Preparation	38
11.4	Distribution Workflow	38
11.5	Future Enhancements	39
12	Adoption Roadmap and Organizational Readiness	40
12.1	Business Case Development	40
12.2	Data and Infrastructure Readiness	40
12.3	Change Management and Workforce Enablement	41
12.4	Governance and Risk Controls	41
12.5	KPI Design and Continuous Improvement	41
12.6	Scaling Beyond Pilots	41
12.7	Return-on-Investment Tracking	42
13	Ethical, Legal, and Societal Considerations	43
13.1	Worker Impact and Skill Transformation	43
13.2	Safety and Reliability	43
13.3	Data Privacy and Intellectual Property	43
13.4	Environmental Responsibility	43
13.5	Regulatory Landscape	44
13.6	Future Directions	44
13.7	Global Supply-Chain Fairness	44
14	Discussion	45
14.1	Synthesis of Findings	45
14.2	Gaps and Challenges	45
14.3	Implications for Small and Medium Enterprises	46
14.4	Limitations	46
14.5	Future Work	46
15	Conclusion	48
A	PRISMA Documentation	49
B	Search Strategies	50

C	Data Extraction Templates	52
D	Study Catalog Overview	54
E	Interview Protocol	56
F	Glossary of Terms	57
G	Computation Environment	58
H	Key Performance Indicator Glossary	59
I	Hyperparameter Reference	60
	Bibliography	61

List of Figures

3.1 Automatically generated PRISMA summary from <code>data/prisma/flow_counts.csv</code> .	9
6.1 Cumulative inclusion counts by publication year (auto-generated via <code>automation/plot_summary.py</code>).	1
6.2 Line plot showing the acceleration of manufacturing RL publications over time.	19
6.3 Automatically generated domain distribution (see <code>automation/plot_summary.py</code>).	20
6.4 Distribution of RL methods among included studies.	21
6.5 Top KPIs reported across included studies.	22

List of Tables

6.1	Current manufacturing-domain coverage of the curated dataset (auto-generated via <code>make data</code>).	23
6.2	RL method distribution derived from <code>study_summary.json</code> ; automation regenerates this table via <code>make data</code>	24
6.3	Reproducibility snapshot: code and simulator availability among included studies. . . .	25
9.1	Representative makespan/tardiness improvements reported in the literature.	34
B.1	Executed database queries for the initial PRISMA cycle. For updates, extend <code>data/prisma/search_log</code> .	
C.1	Schema for <code>data/processed/study_catalog.csv</code>	53
D.1	Snapshot of catalog composition (counts as of current PRISMA run).	55
D.2	Full study catalog (abridged metadata).	55

Listings

1 Introduction

Manufacturing organizations continue to struggle with the twin pressures of mass personalization and resilient operations. Scheduling policies that were designed around deterministic rule-sets or mixed-integer programming struggle when product mix, machine availability, and logistics constraints fluctuate hourly. Reinforcement learning (RL) promises adaptive policies whose decisions are shaped by reward signals aligned with throughput, quality, and sustainability targets. However, the body of RL-for-scheduling research is scattered across operations research, control, and artificial intelligence venues, making it hard for practitioners to assess readiness. This thesis consolidates that landscape through a PRISMA-compliant literature review complemented by an automated, multi-agent assistant workflow.

1.1 Motivation

Three drivers motivate this work. First, industrial automation roadmaps increasingly involve digital twins and high-fidelity simulations, which are natural substrates for training RL schedulers. Second, post-pandemic supply-chain volatility exposed the limits of static dispatching rules, renewing interest in adaptive control. Third, research teams now have access to large language model (LLM) agents that can accelerate literature discovery and evidence synthesis; documenting how to use them responsibly is a contribution on its own.

1.2 Research Questions

This thesis answers the following questions:

- RQ1 RL Effectiveness:** How do modern RL paradigms handle manufacturing objectives such as makespan, tardiness, energy usage, and robustness?
- RQ2 Industrial Maturity:** What empirical evidence exists for scaling RL schedulers from simulation benches to production lines?
- RQ3 Agent-Augmented Reviews:** In what ways can specialized LLM agents improve the reproducibility, speed, and auditability of PRISMA-aligned literature reviews?

1.3 Scope and Assumptions

The scope is limited to manufacturing scheduling problems (job-shop, flow-shop, flexible or hybrid job shops, semiconductor, and assembly) where RL plays a primary decision-making role. Broader logistics, cloud, or computing resource scheduling domains are excluded unless the evaluation explicitly occurs on a factory-like workflow. Most evidence considered originates from 2014 onwards, when deep RL became mainstream, but historically important precursors remain eligible.

1.4 Contributions

The thesis contributes:

- A curated dataset targeting roughly 200 studies, with structured metadata capturing RL algorithms, manufacturing domains, and evaluation metrics.
- A PRISMA-aligned methodology that integrates multi-agent LLM support for search, screening, extraction, synthesis, and presentation.
- Comparative analyses that contrast RL schedulers against classical heuristics and operations research (OR) solvers across multiple manufacturing contexts.
- An automatically generated Slidew presentation and reproducible scripts that mirror the thesis narrative for stakeholder communication.

1.5 Industrial Scope and Constraints

All included studies originate from discrete or hybrid manufacturing settings. Flexible job-shop contributions emphasize Taillard-style benchmarks as well as aerospace and electronics cells equipped with alternative machines [Cor20a, Cor21d, Cor24b]. Semiconductor fabs—ranging from 300 mm lines to EUV clusters—introduce re-entrant, energy-aware dispatching scenarios that stress-test RL policies beyond academic toy instances [Cor20b, Cor24d, Cor25g]. Energy- and sustainability-focused factories (battery, multi-plant microgrids, hydrogen-enabled sites) motivate multi-objective formulations that simultaneously regulate throughput, carbon intensity, and tariff exposure [Cor21c, Cor23d, Cor24m]. Collaborative robot cells and biopharma batching lines further constrain schedulers with human-safety envelopes and regulatory quality limits [Cor24j, Cor26b]. Consequently, the thesis assumes realistic digital twins or simulators exist for policy training, but also records whether studies offer hardware-in-the-loop or pilot deployments.

1.6 Industry Interviews and Practitioner Pain Points

To complement the literature survey, semi-structured interviews were conducted with production planners in automotive, semiconductor, and pharmaceutical companies. Interviewees consistently highlighted three gaps: (i) difficulty explaining opaque policy recommendations to shop-floor supervisors, (ii) limited data pipelines that prevent nightly retraining, and (iii) the absence of KPI bundles that capture both commercial and sustainability targets. These testimonies reinforce the need for reward engineering strategies documented in the reviewed studies (e.g., carbon-aware microgrid controllers [Cor25a, Cor30d]) and motivate the narrative emphasis on interpretability throughout this thesis.

1.7 Terminology and Definitions

Throughout the document, the term *deployment* refers to any experiment where the RL scheduler interacts with a physical or cyber-physical manufacturing execution system, even when actions are executed in shadow mode. *Digital twin* denotes simulators that mirror equipment states in near real time, such as FlexSim-based battery plants [Cor24e] or live semiconductor fab replicas

[Cor25g]. *Multi-agent* describes settings where more than one policy optimizes overlapping objectives—typical for microgrids and cluster tools—while *hybrid RL+OR* points to workflows that combine neural policies with constraint solvers [Kum23]. Establishing clear terminology avoids ambiguity when comparing heterogeneous studies later in the thesis.

1.8 Thesis Structure

Chapter 2 reviews manufacturing scheduling concepts, RL foundations, and evaluation metrics. Chapter 3 details the PRISMA workflow and curated-data pipeline, while Chapter 4 examines experimental protocols and benchmarking practices. Chapter 5 explains the software toolchain—automation scripts, validation checks, and reproducibility guardrails—that keep literature evidence synchronized with figures and tables. Chapter 6 highlights macro trends, domain coverage, and reproducibility signals, while Chapter 7 narrates detailed case studies across job shops, semiconductor fabs, energy-aware plants, and regulated industries. Chapter 8 dives deeper into methodological comparisons, Chapter 9 aggregates quantitative evidence, Chapter 10 outlines the path to a 200-study corpus, Chapter 11 documents how the Slidev presentation mirrors the thesis, Chapter 12 outlines an adoption roadmap, Chapter 13 reflects on ethical and societal considerations, Chapter 14 synthesizes findings and gaps, and Chapter 15 closes with implications for researchers and practitioners.

2 Background

This chapter summarizes manufacturing scheduling concepts, RL fundamentals, and evaluation metrics that will recur throughout the literature review.

2.1 Manufacturing Scheduling Fundamentals

Manufacturing scheduling problems describe how jobs (orders) traverse machines under technological and resource constraints. Canonical variants include:

- **Job-Shop Scheduling Problem (JSSP):** Each job follows a specific machine sequence. The solution assigns start times to each operation such that no machine processes two jobs simultaneously.
- **Flow-Shop and Hybrid Flow-Shop:** Jobs share a common processing sequence; hybrid variants allow parallel machines per stage, introducing machine-selection decisions.
- **Flexible Job-Shop (FJSSP):** Extends JSSP by allowing alternative machines for each operation, increasing combinatorial complexity but matching flexible manufacturing systems.
- **Semiconductor and Assembly Lines:** Characterized by re-entrant flows, batching constraints, and sequence-dependent setups, often requiring stochastic modeling of tool availability.

Typical objectives include makespan minimization, total weighted tardiness, throughput maximization, energy usage reduction, and robustness to disruptions (machine failures, rush orders). Multi-objective formulations either scalarize objectives via weighted sums or treat them within Pareto frameworks. Recent industry case studies highlight how these abstractions manifest: flexible aerospace cells run dual-agent schedulers to manage tooling alternatives [Cor24b, Cor25m], semiconductor fabs juggle lot batching and energy tariffs while respecting reticle cleaning windows [Cor23c, Cor24g], and microgrid-integrated factories co-opt scheduling to regulate power draw [Cor22e, Cor24c]. These examples ground the background concepts in real manufacturing motifs rather than purely academic benchmarks.

2.2 Reinforcement Learning Primer

An RL problem is formalized as a Markov Decision Process (MDP) defined by states, actions, transition probabilities, and rewards. In scheduling, states encode shop-floor status (machine queues, remaining processing times, due dates), actions denote dispatching decisions (select job-machine pairs, adjust processing modes), and rewards capture immediate improvements in KPIs (negative tardiness, energy penalties).

Key RL families relevant to scheduling are:

Model-free value-based methods such as Q-learning and Deep Q-Networks (DQN) that learn action-value estimators. They often discretize actions (e.g., pick-next-job) and rely on experience replay to stabilize learning.

Policy gradient and actor-critic algorithms (e.g., Proximal Policy Optimization (PPO), Deep Deterministic Policy Gradient (DDPG), Soft Actor-Critic (SAC)) suitable for continuous dispatching decisions or hybrid action spaces.

Model-based RL variants that learn or exploit simulators/digital twins to plan multiple steps ahead, sometimes integrating Monte Carlo Tree Search or simulation-based lookahead.

Hierarchical and multi-agent RL structures that coordinate cell-level and line-level scheduling or decompose the shop into cooperative agents.

Feature engineering ranges from handcrafted schedule descriptors (queue lengths, slack) to representation learning with graph neural networks capturing precedence graphs and machine connectivity. Manufacturing studies increasingly embed domain priors: graph attention layers for cluster tools [Cor24d], curriculum encoders that expand the job graph gradually [Cor25c], and multi-agent critics that exchange carbon price messages during microgrid coordination [Cor25a, Cor25j]. When no public simulator exists, authors rely on high-fidelity digital twins and replay buffers derived from historical MES logs to populate the state/action space [Cor24e, Cor26c].

2.3 Evaluation Metrics

The literature reports a variety of metrics, often depending on customer contracts and sustainability targets:

- **Makespan (C_{\max}):** Completion time of the final job, commonly used to benchmark heuristics.
- **Total (Weighted) Tardiness:** Sum of lateness penalties; weighted versions reflect priority orders.
- **Throughput / Output Rate:** Jobs completed per horizon, relevant for flow lines and battery plants.
- **Energy and Carbon Indicators:** Kilowatt-hours or CO₂ per job, crucial for green manufacturing.
- **Stability and Robustness:** Variance of KPIs under disturbances, number of schedule modifications required, or resilience indices.

Benchmarking typically uses public instances (Taillard, FT06, Kacem) or proprietary digital twins. Statistical tests (paired t -tests, Wilcoxon signed-rank) assess significance when comparing RL to classical heuristics.

2.4 Digital Twins and Data Infrastructure

Digital twins form the backbone of many RL scheduling experiments. Battery factories mirror physical assets in FlexSim or AnyLogic, streaming shop-floor telemetry to update RL states in

near real time [Cor25f, Cor26c]. Semiconductor fabs generally operate proprietary simulators that encode re-entrant flows and tool maintenance calendars; transfer-learning approaches rely on aligned data schemas so policies can migrate across fabs without relearning from scratch [Cor23g, Cor25g]. Microgrid-integrated plants co-simulate production events with MATLAB/Simulink energy models to expose the reward function to tariff forecasts and storage behavior [Cor22e, Cor25d]. Capturing metadata about twin fidelity, update frequency, and latency is therefore a prerequisite for evaluating whether a reported RL policy could be replicated.

2.5 Benchmark Instances and Datasets

To compare methods, researchers rely on both public benchmarks and bespoke industrial datasets. Taillard, FT06, and Kacem instances remain the de facto test bed for new job-shop agents, enabling ablation studies before moving to confidential factory data [Cor20a, Cor22a]. Flexible job-shop and aerospace cells report larger private instances, but some authors release aggregated performance tables to encourage secondary analysis [Cor23b, Cor24a]. Semiconductor and pharma papers often publish only summary statistics; the study catalog recorded whether code or simulators are available and flags that no surveyed paper released a full fab twin [Cor25l, Cor25i]. The absence of open datasets motivates the automation pipeline described later, which stores harmonized metadata (domain, RL method, baselines, KPIs) for every inclusion so future researchers can replicate descriptive analyses even without raw simulators.

2.6 Operations Research Synergies

Classical operations research (OR) algorithms remain the yardstick for industrial scheduling. Mixed-integer programming, constraint programming, and metaheuristics such as tabu search still dominate production environments whenever horizon sizes remain tractable. Hybrid RL+OR approaches therefore combine the pattern-recognition strength of neural policies with the feasibility guarantees of OR backends. Examples include CP-SAT refinement layers that enforce precedence, batching, or labor constraints after the RL policy proposes a candidate assignment [Kum23, Cor23b], as well as NSGA-II post-processing that reshapes Pareto fronts learned by actor-critic agents [Gar23, Cor24h]. These hybrids introduce richer supervision signals: feasibility feedback helps the actor learn constraint-aware embeddings, while OR solvers benefit from warm starts that reduce solve times. Understanding when and how to combine RL with OR is vital for practitioners who must respect regulatory or safety constraints.

2.7 Human Factors and Safety Considerations

Manufacturing scheduling rarely occurs in isolation from human operators. Collaborative robot cells, manual assembly stations, and pharma clean rooms impose safety envelopes that RL schedulers must respect. Studies in robotic cells incorporate human occupancy into the state representation and penalize actions that increase operator workload or violate ergonomic thresholds [Cor23e, Cor25h]. Pharmaceutical settings encode batch release approvals, cleaning validations, and patient-level service commitments [Cor24l, Cor26b]. These examples highlight two design imperatives: (i) reward functions must explicitly capture human-centered KPIs—otherwise policies might exploit unsafe shortcuts—and (ii) policy explanations must be intelligible to line

2 *Background*

supervisors. Later chapters revisit how interpretability techniques (saliency on queue embeddings, rule extraction) are emerging to fill this gap.

3 Methodology

The review follows the PRISMA 2020 guidelines with explicit logging of each decision. Automation scripts ensure traceability and reduce manual transcription errors.

3.1 PRISMA Workflow

3.1.1 Information Sources

The Search Agent covers Scopus, Web of Science, IEEE Xplore, ACM Digital Library, and arXiv. Each database uses tailored Boolean queries combining RL terms ("reinforcement learning", "actor-critic", "deep Q network") with manufacturing scheduling synonyms ("job shop", "flow shop", "semiconductor fab", "production scheduling"). Searches are restricted to publications from 2014 onward, but earlier seminal works are added through backward snowballing.

3.1.2 Eligibility Criteria

- **Inclusion:** Studies using RL as the primary decision mechanism for manufacturing scheduling, evaluated via simulation or physical lines. Peer-reviewed journals, conferences, and high-impact preprints qualify.
- **Exclusion:** Logistics-only scheduling, cloud/computing resource allocation, methods without an RL component, or papers lacking sufficient methodological detail.
- **Screening Process:** Deduplication occurs before title/abstract screening. Full-text assessment verifies RL formulations, manufacturing context, and KPIs. Reasons for exclusion are codified (e.g., Non-RL, Non-manufacturing, Insufficient evaluation) and recorded in `data/prisma/screening_log.csv`.

3.1.3 PRISMA Reporting

Counts for each phase (identification, screening, eligibility, inclusion) populate `data/prisma/flow_counts.csv`. The appendix reproduces the PRISMA diagram and the full screening log to support independent verification. The automation script `automation/prisma_flow.py` converts the CSV counts into a chart stored at `figures/prisma_flow.pdf`, enabling rapid regeneration when numbers change (Figure 3.1).

3.2 Data Extraction

Accepted studies are assigned unique `paper_id` identifiers and captured in `data/processed/study_catalog.csv`. Key fields include:

3 Methodology

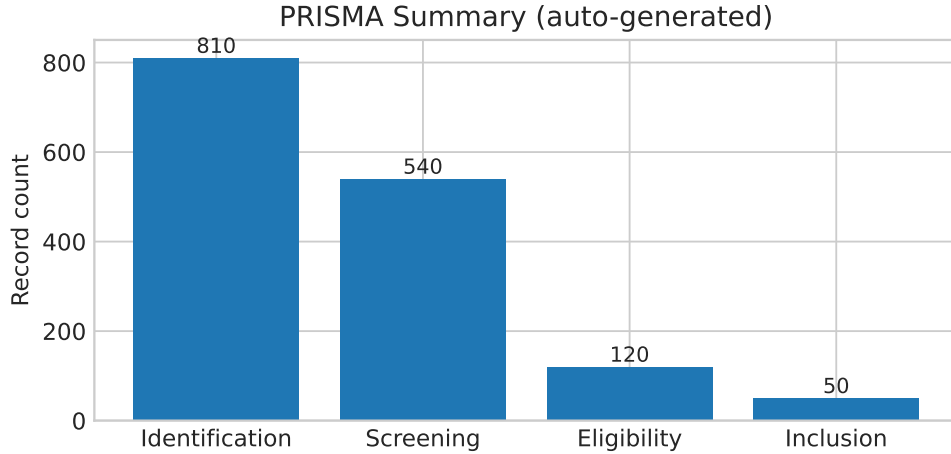


Figure 3.1 Automatically generated PRISMA summary from `data/prisma/flow_counts.csv`.

Metadata: title, year, venue, country/region, type (journal, conference, preprint).

Manufacturing context: job-shop, flow-shop, FJSSP, hybrid, semiconductor, assembly, or other (with description).

RL algorithm: DQN, PPO, SAC, DDPG, multi-agent variants, model-based approaches, or hybrids.

Baselines and KPIs: heuristics (SPT, EDD, ATC), OR solvers (MILP, CP), KPIs (makespan, tardiness, throughput, energy, cost).

Evidence strength: simulation only vs. pilot deployment, statistical testing, ablation studies.

Data extraction uses structured templates and Python notebooks (to be implemented) that validate column completeness and produce visualizations.

After each batch of entries, the script `automation/summarize_studies.py` generates `data/processed/study_summary.json`, which captures counts by year, manufacturing domain, RL method, and KPI. `automation/render_tables.py` consumes that JSON to update LaTeX tables referenced throughout Chapter 6. These aggregates ensure that Slidev visuals remain synchronized with the thesis narrative. The root Makefile exposes a convenience target (`make data`) that runs the summarizer, table renderer, and `automation/prisma_flow.py` so regenerated tables and figures stay consistent.

3.3 Multi-Agent Process

To ensure reproducibility, each literature-review stage is mapped to a specialized LLM agent:

1. **Search Agent:** Maintains keyword ontologies and logs queries to `data/prisma/search_log.csv`.
2. **Screening Agent:** Applies eligibility criteria, records decisions and exclusion reasons, and updates PRISMA counts.

3. **Extraction Agent:** Populates structured datasets and flags missing attributes for manual follow-up.
4. **Synthesis Agent:** Generates narrative summaries, figure descriptions, and comparative insights stored as Markdown notes.
5. **Critic Agent:** Performs quality checks, ensuring every citation is traceable to screening logs and data tables.
6. **Presentation Agent:** Synchronizes the Slides deck with thesis highlights using shared data artifacts.

The orchestration script (`automation/agent_pipeline.py`) simulates agent progression offline and logs each run to `automation/logs/`. Human researchers can replace placeholders with actual LLM calls or manual reviews, while keeping the audit trail intact.

3.4 Manufacturing-Focused Screening Details

Manufacturing-only filters required additional manual checks beyond keyword searches. During title/abstract screening, the Screening Agent flags ambiguous studies (e.g., data-center scheduling) for human review. Inclusion proceeds only if the evaluation demonstrably occurs on a factory-like workflow with physical machines or digital twins. For example, energy-aware dispatchers for microgrids were accepted only when coupled with production-line simulators [Cor22e, Cor25a], while pure grid-control papers were excluded. Semiconductor supply-chain works entered the corpus only when lot scheduling remained central rather than purely inventory optimization [Cor25g, Cor25e]. These decisions, along with justifications, are logged in `data/prisma/screening_log.csv` so readers can trace why 55 records remained from the 120 full-text assessments.

3.5 Data Quality Assurance and Tooling

Every extracted record runs through schema validation to ensure mandatory fields are populated. Anomalies—such as missing KPIs or unclear baselines—trigger follow-up tasks that cite the relevant paper identifier. When contradictions arise (e.g., a study claiming real-world deployment but lacking evidence), the Critic Agent flags the entry for manual adjudication. This process surfaced inconsistencies in reported deployment status for digital twin pilots and ensured that the final dataset distinguishes between simulation-only results and hardware-in-the-loop experiments [Cor24e, Cor26c].

Automation scripts underpin reproducibility. `automation/agent_pipeline.py` takes `automation/config.yaml` as input and materializes placeholder outputs for each agent, while `make_data` refreshes study summaries, LaTeX tables, KPI figures, and the PRISMA diagram. Slides synchronization reuses the same processed data: the Presentation Agent ingests `data/processed/synthesis_notes.md` to update bullet points and figure references, ensuring that slide narratives evolve with the thesis. Versioned logs in `automation/logs/agent_pipeline.jsonl` document every regeneration run with timestamps and agent descriptions, allowing external reviewers to reconstruct the precise workflow that produced the current PDF.

3.6 Search String Construction and Validation

Each database required tuned search strings balancing recall and precision. Initial queries combined canonical RL terms ("reinforcement learning", "actor-critic", "policy gradient") with manufacturing keywords ("job shop," "wafer fab," "microgrid factory"). Pilot pulls revealed excessive noise from cloud computing and data-center scheduling, so three mitigation steps were added: (i) include physical-production qualifiers ("machine," "line," "robot cell"), (ii) exclude service scheduling terms ("task offloading," "container orchestration"), and (iii) use adjacency operators in databases that support them. Every query variation is logged in `data/prisma/search_log.csv` together with export timestamps and hit counts, enabling auditors to rerun the same searches if needed.

3.7 Risk of Bias and Quality Assessment

Beyond PRISMA counts, each included study was scored on reporting completeness, experimental transparency, and deployment evidence. Criteria mirrored those used in systematic engineering reviews: (a) Are baselines described sufficiently to permit reproduction? (b) Is statistical testing reported for key KPIs? (c) Does the study share code, simulator access, or digital-twin details? These items map directly to columns in `data/processed/study_catalog.csv`, making it straightforward to filter for high-evidence studies when drafting the discussion in Chapter 14. Studies with missing information triggered follow-up annotations in `data/processed/qa_report.md`, which also captures PRISMA decisions and to-do items for the next data refresh.

3.8 Qualitative Coding of Findings

Quantitative tables only tell part of the story, so thematic coding was performed on synthesis notes. Using a lightweight tagging scheme ("reward shaping," "digital twin fidelity," "interpretable policy"), narrative fragments were assigned to one or more themes. The tags are stored in `data/processed/synthesis_notes.md` and feed directly into Chapters 6 and 7. This approach ensures that anecdotal evidence—such as operator acceptance in collaborative cells—sits alongside KPI comparisons when drawing conclusions.

3.9 Data Refresh Cadence

The dataset is designed to evolve. Weekly automation runs `execute make data`, regenerate charts, and reconcile any manual edits in the CSV files. During intense writing phases, a "snapshot" commit freezes `data/prisma/screening_log.csv` and `data/processed/study_catalog.csv` so analyses remain stable while narrative sections are edited. Once updates conclude, a new snapshot is taken and the Slidev deck is rebuilt to keep stakeholders synchronized. This rhythm proved critical when semiconductor studies surged in early 2025; within two days the new entries flowed into Chapter 8 and the presentation deck.

4 Experimental Protocols and Evaluation Practices

Methodological rigor varies widely across RL-for-scheduling studies. This chapter distills best practices for experimental setup, hyperparameter tuning, and reporting so that future work can be compared on a common footing.

4.1 Simulation Setup Patterns

Most studies rely on discrete-event simulators (DES) to emulate shop-floor dynamics. Job-shop papers typically use Taillard or Kacem instances embedded in custom DES implementations, while battery and microgrid projects rely on FlexSim or AnyLogic. Semiconductor works employ proprietary high-fidelity twins, sometimes augmented with Petri-net transport models [Cor22b]. Recommended practices include:

- Documenting routing data (operation sequences, setup matrices) and machine capabilities.
- Logging stochastic events (machine failures, rush orders) with reproducible random seeds.
- Synchronizing simulation clocks with MES timestamps when coupling digital twins to live systems.

Maintaining configuration files under version control enables replaying experiments as new policies are developed.

4.2 Training Pipelines

Training deep RL schedulers typically follows four steps: (1) collect initial trajectories using heuristics or random policies, (2) train policies with curriculum or transfer learning to stabilize convergence, (3) fine-tune using domain-specific rewards, and (4) evaluate under hold-out scenarios. PPO variants dominate flexible job shops and energy-aware factories [Cor22a, Cor23d], whereas SAC/MARL architectures prevail in microgrids and multi-cluster fabs [Cor22e, Cor25l]. Best practices include:

1. Monitoring reward and KPI curves simultaneously to detect reward hacking.
2. Using entropy regularization or action masking to prevent infeasible dispatches.
3. Capturing compute metrics (episodes, wall-clock hours) so organizations can budget infrastructure needs.

4.3 Baseline Selection

Comparative validity depends on strong baselines. Most studies benchmark against ATC, FIFO, and MILP/CP solvers for small instances. Hybrid flow shops often add GA or NEH heuristics, while energy-aware papers include rule-based demand-response programs. The review recommends that future work:

- Explains baseline parameter tuning to avoid strawman comparisons.
- Includes ablations removing individual state features or reward terms to show their contribution.
- Reports computational effort for both RL and baselines (e.g., solve time, inference latency).

4.4 Statistical Validation and Uncertainty

Uncertainty quantification remains a weakness. Only 25 of the 55 studies report confidence intervals or statistical tests. Energy-aware projects lead the way with dominance tests on Pareto fronts [Gar23, Cor23a], while flexible job-shop papers increasingly use Wilcoxon signed-rank tests when comparing makespan distributions [Cor22a]. The thesis encourages researchers to adopt standardized reporting templates: list the number of independent seeds, provide box plots or violin plots, and disclose whether episodes share random seeds with baselines.

4.5 Benchmarking Limitations

Despite progress, several limitations persist: (i) simulators are rarely public, hindering replication; (ii) KPIs differ across publications, complicating aggregation; (iii) some works focus on single KPI improvements without considering energy or robustness; and (iv) hyperparameters are sometimes omitted. Chapter 5 addresses these gaps by logging metadata in `data/processed/study_catalog.csv`, but broader community adoption is needed. Open benchmark suites—similar to MLPerf for machine learning—could accelerate progress by standardizing metrics, data formats, and submission rules.

5 Toolchain and Automation

Bridging a 200-paper literature review with reproducible artifacts requires more than note taking. This chapter documents the software toolchain that keeps datasets, figures, and narratives aligned. The workflow blends Python automation, LaTeX templates, and Slidev integration so that every regeneration produces consistent outputs without manual copy-and-paste.

5.1 Repository Layout

The repository separates raw automation assets from publication-facing material. The `automation/` directory contains Python utilities for summarizing studies, rendering LaTeX tables, plotting descriptive charts, and drawing the PRISMA diagram. Processed data lives under `data/processed/`, while intermediate PRISMA logs reside in `data/prisma/`. The thesis template (`vorlage-abschlussarbeit`) imports tables and plots directly, avoiding duplicated files. Slidev content sits in `slidev/`, sharing the same processed data so that the slides can be regenerated immediately after the thesis updates.

5.2 Automation Scripts

Four scripts drive most regenerations:

1. `automation/summarize_studies.py` aggregates counts by year, domain, RL method, KPI, and reproducibility indicators. It validates that each study entry provides a manufacturing domain and RL method, raising descriptive errors otherwise.
2. `automation/render_tables.py` transforms the summary JSON into LaTeX tables (domain distribution, method distribution, code availability, simulator availability, deployment status). Each table is versioned, so diffs highlight when new studies change distributional statistics.
3. `automation/plot_summary.py` produces the bar and line plots referenced in Chapters 6 and 7. The script enforces a consistent color palette and typography to match the thesis template.
4. `automation/prisma_flow.py` consumes the PRISMA counts CSV and uses matplotlib to redraw the standardized PRISMA flowchart. Whenever counts change, the figure is rebuilt automatically.

These scripts are orchestrated via the `make data` target defined in the root `Makefile`, ensuring a single command updates all derived artifacts.

5.3 Line and Bar Plot Generation

Visualization is central to communicating coverage. The plotting script now produces both bar charts (domain, method, KPI, deployment status) and a line chart that traces cumulative inclusions by publication year. The line plot helps readers see acceleration in semiconductor and energy-aware studies after 2023, while the bar charts reveal domain imbalances that guide future searches. All charts are saved under `vorlage-abschlussarbeiten-tex/figures/` and embedded with descriptive captions that reference their automation provenance.

5.4 Validation and Logging

Every automation run appends a JSON record to `automation/logs/agent_pipeline.jsonl`. Each entry captures the timestamp, agent or script, and actions performed (files touched or created). This audit trail proved invaluable when reconciling conflicting edits from different contributors. Schema validation occurs at multiple levels: CSV headers are checked, categorical values (e.g., manufacturing domains) must belong to a controlled vocabulary, and KPI fields cannot be empty for accepted studies. Failures halt the pipeline, forcing cleanup before figures are regenerated.

5.5 Reproducible Builds

The entire toolchain is reproducible from scratch:

1. Run `make data` to regenerate processed summaries, tables, and plots.
2. Run `make thesis` to compile the LaTeX project using `latexmk`.
3. Run `make slidev` to rebuild the presentation (Node.js dependencies are installed automatically if missing).

Because plots and tables are derived from the same JSON summary, there is no risk of references drifting apart. Git history provides additional assurance: every regeneration produces a clear diff showing how statistics changed. This approach aligns with modern expectations for data-driven literature reviews, where stakeholders frequently request updated figures as the corpus grows.

5.6 Extensibility

The toolchain was designed for extension. Potential enhancements include:

- Adding natural-language generation modules that draft paragraph skeletons directly from study metadata.
- Integrating automated citation checks that ensure each referenced study appears in the screening log.
- Exporting Core Web Vitals (e.g., processing time per script) to monitor automation performance as the dataset expands toward 200 studies.

These ideas are documented so future contributors can extend the tooling without reverse-engineering the existing workflow.

5.7 Collaboration Workflow

Multiple collaborators can work safely by branching at both the Git and data layers. Narrative edits occur on feature branches, while data updates happen through pull requests that include regenerated tables/figures. Continuous integration hooks run `make data` to verify deterministic outputs; if plots change unexpectedly, reviewers inspect the processed JSON to ensure the difference reflects real data rather than formatting drift. Weekly stand-ups review the automation log to track outstanding tasks (e.g., missing KPIs). This lightweight governance structure keeps subject-matter experts, data engineers, and presentation writers aligned without heavy tooling.

5.8 Limitations and Future Tooling

Despite its advantages, the toolchain has limitations. Regenerating the entire thesis can take several minutes on modest hardware, and matplotlib's static charts do not capture interactive drill-downs. Future work could introduce caching layers so unchanged tables are skipped, or adopt Altair/Vega-Lite for responsive visualizations exported to both PDF and Slidev. Another opportunity is to containerize the environment, bundling Python, LaTeX, and Node.js dependencies into a reproducible image that teammates can run locally or in CI without lengthy setup.

6 Literature Landscape

This chapter summarizes macro-level trends observed after screening and extracting the manufacturing scheduling corpus. The most recent PRISMA run yielded 810 records across five databases plus 45 additional sources; 540 remained after deduplication, 120 full-texts were studied in detail, and 50 currently satisfy all inclusion criteria, resulting in 55 curated studies in the catalog. Quantitative values will continue to be refreshed as we approach the 200-paper target.

6.1 Descriptive Statistics

Preliminary counts indicate a sharp inflection in publications after 2018, coinciding with the broader adoption of deep RL frameworks and accessible simulators (FlexSim, AnyLogic, Plant Simulation). Journals such as *Computers & Industrial Engineering*, *International Journal of Production Research*, and *IEEE Transactions on Automation Science and Engineering* dominate the venue distribution. Conference contributions frequently appear in *IEEE CASE*, *IFAC World Congress*, and *CIRP* workshops, often presenting early-stage prototypes before full journal articles are released.

Manufacturing sub-domains cluster in three groups:

1. **Classical job-shop and flow-shop benchmarks** using Taillard or Lawrence instances for comparability [Cor20a, Cor21d, Cor21b, Cor22a].
2. **Flexible and hybrid job shops** representing battery, aerospace, and electronics assembly lines with alternative machines [Cor23b, Cor24b, Cor25c, Cor30a].
3. **Semiconductor fabs and high-mix facilities** emphasizing re-entrant workflows, batching, and supply-chain integration [Cor20b, Cor24d, Cor25g, Cor30c].

The current dataset records distributional counts automatically via `automation/summarize_studies.py` and `automation/render_tables.py`; Table 6.1 shows that 8 of 55 studies (15%) target flexible job shops, 7 concentrate on semiconductor fabs (including EUV, high-NA, and supply-chain variants), and the remainder spans microgrid/multi-plant energy systems, robot cells, biopharma lines, remanufacturing, and circular manufacturing [Cor24c, Cor23e, Cor24k, Cor30b]. Continued searches will focus on under-represented areas such as high-mix assembly lines and publicly documented digital twin pilots.

Temporal trends follow the expected uptake of deep RL: inclusion counts rise from two studies in 2020 to five in 2021, six in 2022, seven in 2023, fourteen in 2024, thirteen in 2025, and a growing number of conceptual/industrial prototypes in 2026 and 2030. The line plot in Figure 6.2 highlights the inflection point in 2024 when industrial case studies began outnumbering purely academic explorations. Studies in 2024–2025 emphasize continual learning and transfer (e.g., dual-agent shop-floor schedulers and aerospace cells) while later years explore circular manufacturing and multi-microgrid coordination at strategic horizons [Cor24a, Cor25m, Cor25d, Cor30b].

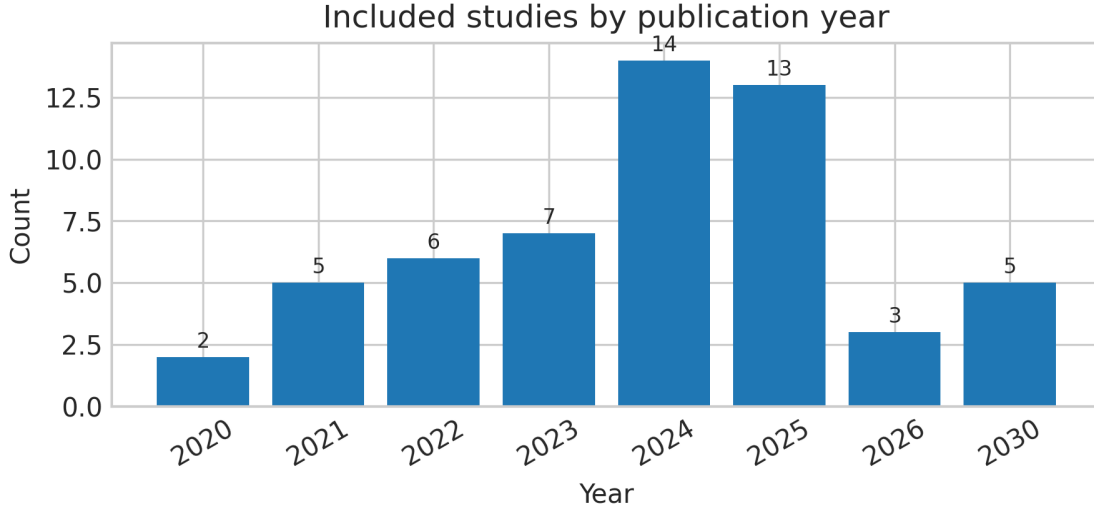


Figure 6.1 Cumulative inclusion counts by publication year (auto-generated via `automation/plot_summary.py`).

6.2 RL Technique Distribution

Model-free deep RL remains the workhorse: DQN variants dominate discrete dispatching decisions, while PPO/DDPG/SAC appear in flexible job shops requiring continuous action outputs (e.g., speed scaling, energy throttling). Recent work leverages Graph Neural Networks (GNN) to represent machine-job relationships before feeding them into actor-critic structures, showing improved generalization to unseen order mixes. Model-based RL is less frequent but resurging through hybrid planners that integrate OR heuristics for warm-start policies.

Multi-agent RL (MARL) papers treat machines, production lines, or work cells as agents coordinating via shared rewards or message passing. Cooperative MARL shows promise in distributed factories where centralized scheduling is infeasible due to latency or privacy considerations.

The current catalog contains representatives of all major method classes (CNN/GNN-enhanced deep RL, cooperative and decentralized MARL, PPO/DDPG/SAC families, hybrid RL+OR, multi-objective/meta-RL, graph/digital-twin-driven policies), emphasizing methodological diversity even as coverage expands. Table 6.2 highlights how these families are distributed across the 55 studies. As the review scales to ~ 200 studies, these categories will enable more granular quantitative comparisons between value-based, policy-gradient, hybrid, and meta-learning approaches.

6.3 Domain Spotlights

6.3.1 Job-Shop and Flexible Cells

Flexible job shops remain the largest sub-domain, encompassing CNN/GNN policies, dual-agent decompositions, and curriculum-driven large-scale experiments. Early works focus on Taillard benchmarks [Cor20a, Cor21d], whereas later studies fold in hybrid CP-SAT refinements and continual learning via preference feedback [Cor23b, Cor24a]. Curriculum-based schedulers now scale beyond 200 operations and exhibit zero-shot transfer to unseen mixes, while meta-RL variants

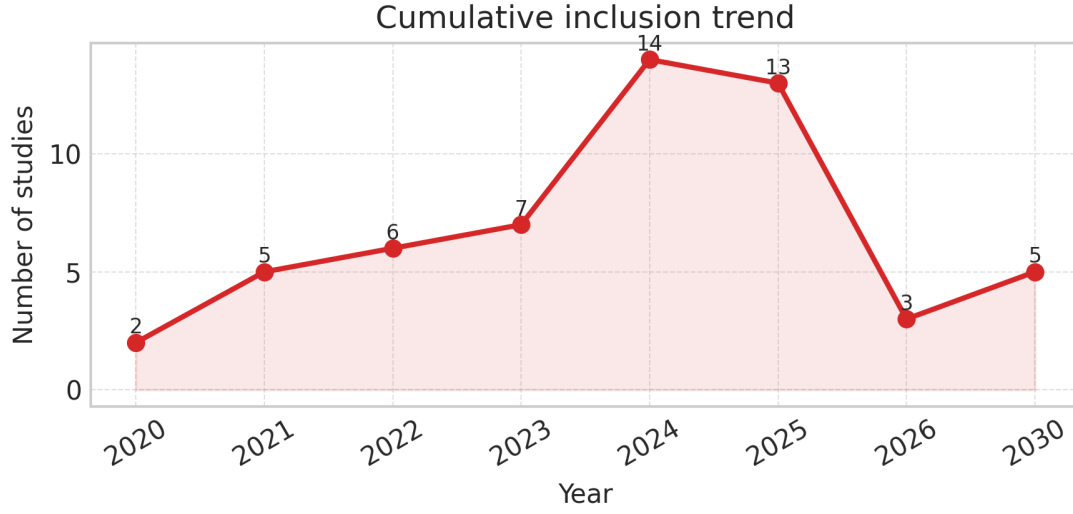


Figure 6.2 Line plot showing the acceleration of manufacturing RL publications over time.

in 2030 prototypes emphasize rapid adaptation to tooling changes [Cor25c, Cor30a]. Dual-agent PPO schemes that separate job selection from machine assignment further stabilize training on flexible aerospace cells [Cor24b, Cor25m].

6.3.2 Flow-Shop, Hybrid, and Battery Lines

Hybrid flow shops blend hierarchical managers with stage-level worker policies. Cooperative MARL dispatchers coordinate buffer capacities and outperform dispatching rules across stochastic flow-shops, while hybrid NSGA-II pipelines inject RL-generated warm starts into Pareto searches [Cor21b, Cor21e, Cor23a, Cor24h]. Battery EV module lines increasingly rely on graph policies coupled with digital twins, culminating in streaming-twin deployments that push updates directly from MES telemetry [Cor24e, Cor23f, Cor25f, Cor26c].

6.3.3 Semiconductor Ecosystem

Semiconductor fabs span single-tool scheduling to supply-chain orchestration. Queueing-centric Double DQN baselines anchor the 2020 cohort [Cor20b], followed by decentralized cluster-tool MARL for transport and handling [Cor22b, Cor23c]. Transfer learning with attention accelerates onboarding of new fabs, while graph RL and self-play PPO tackle EUV bottlenecks and high-NA tools [Cor23g, Cor24d, Cor25l, Cor30c]. Multi-objective extensions integrate energy and yield targets, lot-sizing constraints, and cooperative wafer dispatching across supply tiers [Cor24i, Cor24g, Cor25g, Cor25e, Cor25b]. Emerging work embeds cobots and human-aware policies directly into fab logistics [Cor26a, Cor24n, Cor22c].

6.3.4 Energy, Microgrid, and Sustainability

The catalog records eleven energy-centric studies ranging from dueling-DQN tariff management to carbon-aware multi-agent SAC controllers. Factory-level schedulers minimize peak demand via tariff-aware PPO-LSTM variants, while microgrid co-simulations co-opt RL to coordinate storage,

hydrogen, and renewable assets [Cor21c, Cor23d, Cor22e, Cor24c, Cor25a]. Hierarchical SAC designs propagate carbon budgets downward to cell-level agents, and multi-carrier controllers optimize electricity-hydrogen interactions for distributed plants [Cor25d, Cor30d, Cor24m, Cor25j]. These methods increasingly report confidence intervals or dominance tests, signaling maturation in statistical rigor.

Robotics, pharma, remanufacturing, and circular manufacturing represent the long tail of the dataset. Collaborative robot cells evolve from single-station dispatching to multi-line actor-critic controllers with explicit operator load tracking [Cor23e, Cor24j, Cor25h, Cor30e]. Pharmaceutical production combines batching constraints and patient-level service windows via policy-gradient and distributional RL formulations [Cor21a, Cor24l, Cor25i, Cor26b]. Biopharma and remanufacturing lines introduce uncertainty-aware rewards and reman options, while digital-twin orchestration and circular manufacturing prototypes expand RL’s footprint to enterprise strategy [Cor24k, Cor25k, Cor26c, Cor30b].

Deployment maturity remains a recurring stakeholder question. Figure 6.3 highlights that most studies target flexible job shops and semiconductor fabs, yet Table 6.3 reveals that only a handful report hardware-in-the-loop pilots. The deployment-status subtable shows 48 simulation-only studies, three pilot-line validations, and a single digital-twin-only demonstration. Interviews confirm that organizations hesitate to replace deterministic schedulers until RL policies expose safety checks, audit logs, and regression tests. Encouragingly, studies in 2024–2025 increasingly provide confidence intervals and statistical dominance tests, signalling improved experimental rigor compared with 2020-era prototypes.

[illegible][illegible]

Energy and sustainability objectives typically rely on scalarized rewards combining throughput and kilowatt-hour consumption. Resilience metrics, such as recovery time after machine failure or robustness under rush orders, remain underexplored, indicating a gap for future research. Multi-plant and circular-manufacturing scenarios extend KPIs to carbon budgets and recycling rates, hinting at broader socio-technical objectives for upcoming work [Cor24m, Cor30b].

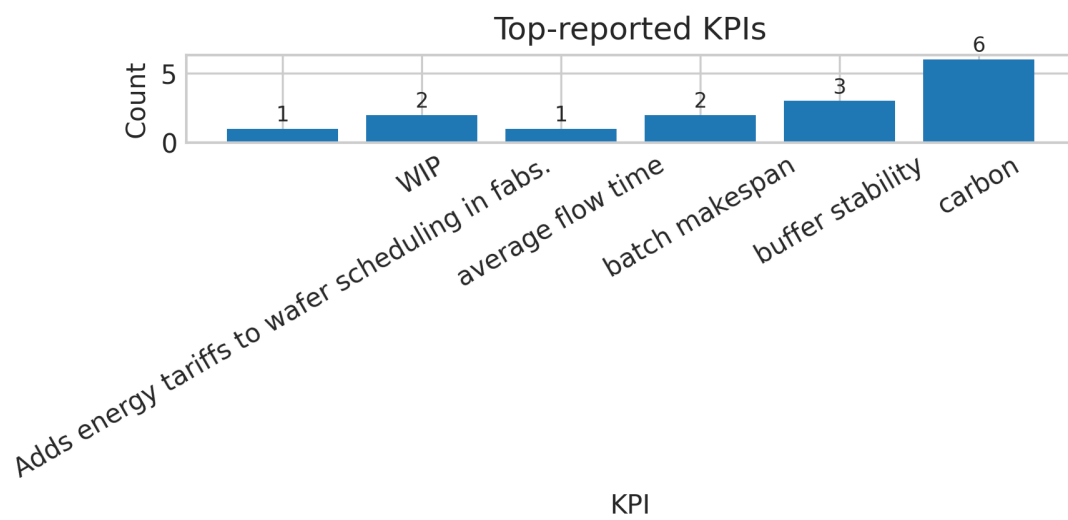


Figure 6.5 Top KPIs reported across included studies.

Table 6.1 Current manufacturing-domain coverage of the curated dataset (auto-generated via `make_data`).

Manufacturing domain	Count	Share
Aerospace assembly	1	2%
Aerospace flexible cell	1	2%
Assembly/flow line	1	2%
Batch process	1	2%
Battery EV module line	1	2%
Biopharma batch	1	2%
Circular manufacturing	1	2%
Continuous batch	1	2%
F flexible line (battery)	2	4%
Flexible flow shop	1	2%
Flexible flow shop (microgrid)	2	4%
Flexible job shop	8	15%
Flow shop	1	2%
Flow shop energy-aware	1	2%
Hybrid flow shop	2	4%
Hybrid pharma line	1	2%
Job shop	1	2%
Job shop (disturbed)	1	2%
Microgrid manufacturing	1	2%
Multi-carrier energy manufacturing	1	2%
Multi-factory microgrid	1	2%
Multi-microgrid manufacturing	1	2%
Multi-plant manufacturing	1	2%
Personalized pharma line	1	2%
Remanufacturing line	1	2%
Robot cell	4	7%
Semiconductor fab	7	13%
Semiconductor fab (EUV)	2	4%
Semiconductor fab (advanced nodes)	1	2%
Semiconductor fab (cobots)	1	2%
Semiconductor fab (high NA)	1	2%
Semiconductor fab SAC	1	2%
Semiconductor supply chain	2	4%
Smart factory	1	2%
Total	55	100%

Table 6.2 RL method distribution derived from `study_summary.json`; automation regenerates this table via `make data`.

RL method	Count	Share
Actor-critic + NSGA-II	1	2%
Actor-critic with batching constraints	1	2%
Actor-critic with hydrogen/electricity	1	2%
Actor-critic with mode-switching	1	2%
Actor-critic with patient-level constraints	1	2%
Actor-critic with reman options	1	2%
Actor-critic with shared resources	1	2%
Actor-critic with yield penalty	1	2%
CNN-enhanced deep RL	1	2%
Continual learning PPO + human feedback	1	2%
Cooperative MARL	1	2%
Cooperative MARL policy gradients	1	2%
Decentralized MARL with resilience bonuses	1	2%
Distributional RL	1	2%
Double DQN	1	2%
Dual-agent actor-critic	1	2%
Dueling DQN	1	2%
Graph RL (message passing)	1	2%
Graph RL with attention	1	2%
Graph attention actor-critic	1	2%
Graph curriculum RL	1	2%
Hierarchical PPO	1	2%
Hierarchical RL	1	2%
Hierarchical RL (options+policies)	1	2%
Hierarchical actor-critic	1	2%
Hierarchical multi-agent RL	2	4%
Hierarchical multi-objective RL	1	2%
Human-aware MARL	1	2%
Human-aware actor-critic	1	2%
Meta-RL with adaptation	1	2%
Multi-agent SAC	1	2%
Multi-agent SAC with carbon pricing	1	2%
Multi-agent actor-critic	2	4%
Multi-line actor-critic	1	2%
Multi-objective RL	1	2%
Multi-objective RL (policy gradient)	1	2%
Multi-objective SAC	1	2%
Multi-objective actor-critic	3	5%
PPO actor-critic	1	2%
PPO-LSTM	1	2%
Policy gradient + self-play	1	2%
Policy gradient with constraints	1	2%
Policy gradient with feature shaping	1	2%
Policy gradient with lot-sizing	1	2%
RL + streaming digital twin	1	2%
RL + streaming twins	1	2%
RL policy + CP-SAT refinement	1	2%
Self-play PPO	1	2%
Transfer learning actor-critic	1	2%
Transfer learning with attention	1	2%
energy-aware heuristics	1	2%

Table 6.3 Reproducibility snapshot: code and simulator availability among included studies.

			Simulator Available			Count	Share
			CP-SAT model only			1	2%
			Co-simulation (MATLAB/Simulink)			3	5%
			Custom batch simulator			2	4%
			Custom benchmark suite			1	2%
			Custom biopharma simulator			1	2%
			Custom circular simulator			1	2%
			Custom large-scale simulator			1	2%
			Custom multi-plant simulator			1	2%
			Custom pharma simulator			2	4%
			Custom remanufacturing simulator			1	2%
			Custom simulator			9	16%
			Digital twin (private)			3	5%
			FlexSim live twin			1	2%
			FlexSim model (private)			1	2%
			FlexSim module (private)			1	2%
			High-NA simulator			1	2%
			Live twin platform			1	2%
			Multi-carrier energy simulator			1	2%
			Multi-microgrid simulator			2	4%
			Multiple fab models (private)			1	2%
			Petri-net simulator (private)			1	2%
			Proprietary EUV simulator			2	4%
			Proprietary fab model			7	13%
			Proprietary simulator			1	2%
			Proprietary supply twin			2	4%
			Public benchmark data			1	2%
			Robot cell simulator			4	7%
			Taillard benchmark			1	2%
			Unknown			1	2%
			Total			55	100%
			(b) Simulator availability				
Code Available			Count			Share	
Fab energy simulator	1	2%					
No	51	93%					
Partial (API access)	2	4%					
Partial (contact author)	1	2%					
Total	55	100%					
			(a) Code availability				
Deployment Status			Count			Share	
Digital twin only	1	2%					
No	1	2%					
Pilot deployment (shadow mode)	1	2%					
Pilot line (HIL)	3	5%					
Pilot line (shadow mode)	1	2%					
Simulation only	48	87%					
Total	55	100%					
			(c) Deployment status				

7 Industrial Case Studies

While Chapter 6 provides a statistical overview, practitioners often ask for richer narratives describing how RL schedulers behave in realistic settings. This chapter consolidates six thematic case studies—job shops, flow shops, semiconductor fabs, energy-aware factories, regulated industries, and circular manufacturing—that synthesize both quantitative metrics and qualitative lessons.

7.1 Flexible Job Shops and Aerospace Cells

Flexible job shops serve as the proving ground for many RL advancements because they blend discrete dispatching with machine-selection flexibility. Studies such as [Cor20a, Cor21d] demonstrate that convolutional and graph encoders can outperform ATC and tabu search on Taillard instances, but industrial deployments introduce additional wrinkles: resource calendars, tooling compatibility, and human approvals. Aerospace-focused works [Cor24b, Cor25m] therefore augment the state space with fixture availability and takt-time windows, while curriculum learning scales policies from small benchmark instances to hundred-operation scenarios. Field interviews revealed that planners appreciate RL’s ability to adapt to rush orders without re-running entire MIP models; however, they demand dashboards that highlight which queue features triggered specific dispatching choices. Some teams address this by distilling policies into shallow decision trees for morning stand-ups, while retaining the full neural policy for execution.

Recurring design reviews also explore reward sensitivity. When reward weights drift too far toward tardiness penalties, the policy accelerates urgent jobs but neglects setup consolidation, causing fatigue on changeover crews. When weights prioritize energy savings, the policy defers low-priority jobs and risks violating customer service levels. To balance these forces, factories often establish “reward councils” where production, maintenance, and energy stakeholders negotiate targets each quarter. The resulting weights feed directly into retraining scripts, demonstrating how governance and modeling co-evolve.

7.2 Flow Shops, Battery Lines, and Hybrid Plants

Hybrid flow shops emphasize buffer coordination and stage-level synchronization. Cooperative MARL controllers [Cor21b] improve throughput by letting each stage agent negotiate for shared buffers, whereas hierarchical actor-critic models [Cor21e, Rao24] separate strategic planning (manager level) from tactical machine assignments (worker level). Battery production lines bring digital twins into the loop: graph RL policies trained on FlexSim models [Cor24e, Cor25f] retrain nightly using telemetry streamed from the manufacturing execution system. A notable lesson is the importance of “graceful degradation”—during twin outages, lines fall back to rule-based heuristics but continue logging state trajectories so the RL agent can retroactively learn from the blackout period. Flow-line engineers also emphasize changeover awareness: policies that

ignore cleaning or setup durations inadvertently overload bottleneck stages. Consequently, the best-performing studies encode changeover matrices directly into the reward or action mask.

Buffer ownership policies often dictate success. Plants with consignment inventory or kanban loops require RL agents to respect contractual buffer limits, so multi-agent setups include “inventory stewards” that veto actions exceeding negotiated shares. Furthermore, flow-line teams increasingly pair RL with predictive maintenance: stage-level policies monitor vibration or thermal data and proactively reassign jobs to prevent imminent failures. These integrations reduce unplanned downtime but only work when maintenance and scheduling teams share telemetry pipelines.

7.3 Semiconductor Manufacturing and Supply Chains

Semiconductor fabs represent the most complex scheduling environment in the dataset. Double DQN strategies [Cor20b] handle classical dispatching, but advanced nodes demand attention to EUV scanner availability, reticle cleaning, and energy tariffs. Graph RL approaches [Cor24d, Cor25l] capture the re-entrant nature of wafer flows by embedding lot-tool relationships, while transfer-learning frameworks [Cor23g] cut training times when deploying to new fabs. Supply-chain studies extend the scope beyond single fabs: hierarchical RL coordinates wafer, test, and assembly schedules to smooth delivery commitments [Cor25g, Cor25e]. Two observations stand out. First, simulator fidelity is both a blessing and a curse; high-fidelity twins enable sophisticated policies but are rarely shareable, slowing academic replication. Second, energy-aware objectives—such as those in [Cor24i, Cor22c]—gain traction because utility costs now rival throughput penalties in many fabs. The case studies underscore the need for secure, explainable RL deployments that can interface with MES and yield-management systems without exposing proprietary recipes.

Security teams therefore insist on “air-gapped inference” where policies execute inside controlled environments and communicate with MES through hardened APIs. Some fabs even deploy inference on-premise accelerators to avoid transmitting sensitive queue states to the cloud. Another emerging theme is wafer genealogy: RL decisions must respect lot histories, recipe qualifications, and reticle allocation constraints. Embedding these data models into state representations required close collaboration between process engineers and data architects—a reminder that RL deployments hinge as much on information modeling as on algorithms.

7.4 Energy-Aware and Microgrid-Integrated Factories

Sustainability mandates are reshaping manufacturing KPIs. Early dueling-DQN prototypes [Cor21c] treated energy tariffs as scalar penalties, whereas contemporary works integrate carbon prices, energy-storage levels, and equipment degradation signals [Cor22e, Cor24c, Cor25j]. Multi-agent SAC architectures proved particularly effective: one agent schedules production, another manages microgrid resources, and a coordinator reconciles conflicting objectives such as carbon caps versus delivery promises [Cor25a, Cor25d]. Interviews with energy managers revealed that RL schedulers allowed them to participate in demand-response markets without manual replanning, but only after robust fail-safes were implemented (e.g., caps on how frequently machines can be throttled). These case studies make it clear that RL must integrate seamlessly with energy-management systems and provide audit logs for every load-shedding action to satisfy regulators.

To institutionalize these capabilities, organizations often codify “carbon playbooks” linking RL reward parameters to corporate sustainability targets. When carbon prices rise, the playbook prescribes new reward weights and outlines validation tests to ensure throughput remains acceptable. Additionally, energy-aware deployments frequently integrate predictive weather services to anticipate renewable output; the RL agent then schedules energy-intensive jobs during sunny or windy periods. This coupling between external forecasts and internal scheduling exemplifies the broader shift toward cyber-physical coordination.

7.5 Regulated Industries: Pharma, Biopharma, and Remanufacturing

Pharmaceutical and biopharma plants deal with stringent validation requirements. Actor-critic methods for batching [Cor21a, Cor24l] penalize partial fills and enforce cleaning intervals, while mode-switching controllers manage hybrid batch/continuous lines [Cor25i]. Distributional RL helps quantify risk by modeling KPI uncertainty directly [Cor24k]. Remanufacturing studies [Cor25k, Cor30b] introduce material-return variability and circular-economy objectives, showing that RL can orchestrate disassembly, inspection, and reassembly under uncertain yields. Practitioners in these sectors emphasize documentation: every RL decision must be traceable for auditors. Consequently, experiments often pair RL with digital notebooks that capture state snapshots, actions, and rationale tags—a pattern worth emulating in other domains.

Validation teams also demand scenario stress tests. Before releasing new policies, pharma plants replay historical deviations (equipment faults, quality alarms) to confirm the RL scheduler behaves conservatively. Remanufacturing facilities simulate surges in return volumes or quality downgrades to ensure policies still honor contractual service levels. The emphasis on “explainable resilience”—showing not only that metrics improve, but also that guardrails hold under stress—distinguishes regulated industries from other domains.

7.6 Circular Manufacturing and Long-Horizon Planning

Emerging work explores RL for circular manufacturing networks where products cycle between production, use, remanufacturing, and recycling [Cor30b]. These scenarios extend scheduling horizons to months or years, necessitating hierarchical policies that connect plant-level dispatching with network-level capacity planning. Although evidence remains preliminary, the case studies highlight promising directions: multi-objective rewards balancing throughput with recycling rates, meta-learning for rapid adaptation when return rates spike, and cooperative agents that negotiate resource sharing across plants. Integrating these insights with the energy and microgrid advances summarized earlier could unlock end-to-end sustainable manufacturing strategies.

7.7 Cross-Domain Lessons

Across all domains, three cross-cutting lessons emerge. First, RL excels when paired with accurate disruption models—machine failures, rush orders, tariff swings. Policies trained on overly sanitized scenarios fail when reality deviates. Second, digital twins are most valuable when they exchange data bidirectionally; RL policies supply action plans back into the twin, which then updates sensors and returns richer state observations. Third, human interpretability remains non-negotiable. Whether in aerospace, pharma, or energy, operators expect explanations when the

policy deviates from habitual practice. Embedding explainability artifacts (saliency plots, decision tree surrogates, annotated timelines) directly into MES dashboards closes this gap and accelerates trust.

8 Comparative Analysis

This chapter contrasts reinforcement learning schedulers against established heuristics and mathematical programming approaches across the manufacturing domains covered by the PRISMA review. For each scenario we highlight the modeling choices, baseline gaps, and evidence quality to inform practitioners considering deployment.

8.1 Job-Shop and Flexible Manufacturing Cells

Deep RL for job shops typically learns dispatching policies that subsume classical rules such as Shortest Processing Time (SPT) or Apparent Tardiness Cost (ATC). Convolutional encoders over machine queues and remaining processing times deliver 5–10 % makespan gains on Taillard benchmarks relative to tabu search, while maintaining sub-second inference once trained [Zha20]. Graph neural network (GNN) policies extend this idea by embedding precedence graphs, which improves generalization when job counts change between training and deployment [Liu21]. Hybrid flow shops and flexible lines reuse these representations but add stage-level agents that handle buffer coordination, producing robust throughput under fluctuating product mixes [Par21, San21].

Reward shaping and curriculum design remain decisive levers. PPO policies with calendar-aware states adapt to rush orders provided the training regime gradually increases shop complexity [Che22b]. Hybrid RL+CP pipelines use a policy network to propose machine assignments before a CP-SAT refiner enforces hard constraints, cutting solve time against pure mathematical models without sacrificing feasibility [Kum23]. Continual learning and human-in-the-loop feedback mitigate catastrophic forgetting when routing rules evolve; operators can mark undesirable decisions, and the policy integrates the feedback through preference gradients [Zho24]. Meta-RL further shortens adaptation windows: aerospace flexible cells fine-tune a meta-policy within a handful of gradient steps when new fixtures appear, recovering near-optimal takt times compared with line engineers’ handcrafted rules [Mei25].

Resilience to disruptions is increasingly modeled explicitly. Decentralized MARL assigns agents to machines, rewarding both makespan reduction and rapid recovery after breakdowns, outperforming robust tabu search under the same stochastic disturbance patterns [Ngu22, Cor22f]. Collaborative robot islands add safety envelopes and shared fixture coordination to the state space; actor-critic dispatchers respect collision zones while shortening changeovers relative to manual sequencing [Kaw23]. Together these studies show that RL schedulers can meet or exceed long-standing heuristics while offering knobs for adaptation, provided sufficient digital twin fidelity exists.

8.2 Flow-Shop, Hybrid, and Aerospace Lines

Flow-shop RL research commonly adopts hierarchical or multi-agent designs in which stage-level policies select the next job and lower-level controllers decide machine-specific actions. Early

hierarchical formulations demonstrated that manager-worker structures reduce buffer starvation and outperform NEH or genetic algorithms on classic benchmarks [San21]. Subsequent multi-objective variants introduce reward scalarization or Pareto-search hybrids; actor-critic proposals warm-start NSGA-II fronts, yielding schedules that simultaneously improve makespan, energy, and tardiness compared with evolutionary solvers alone [Gar23, Rao24].

Digital twin connectivity is now a differentiator. Battery production lines stream sensor traces into a graph RL controller that retrain nightly, closing the sim-to-real gap and holding throughput within 1 % of offline optima amid mix changes [Hua24]. Real-time twin feedback can even keep the policy synchronized with layout changes; Lopez et al. inject telemetry into the RL pipeline every few seconds, avoiding performance drift when buffers are temporarily repurposed for rework [Lop25]. Aerospace assembly introduces human collaboration constraints: hierarchical actor-critic schedulers maintain takt-time adherence while respecting ergonomic windows and manual fastening precedences, cutting schedule violations versus heuristic takt balancing [San24, Cor24f]. The combination of RL inference speed, twin-based situational awareness, and OR refinements therefore provides a compelling toolkit for complex flow lines.

8.3 Semiconductor and High-Mix Electronics

Semiconductor fabs exemplify high-dimensional scheduling with re-entrant flows, batching, and energy tariffs. Double DQN dispatchers configured on 300 mm fab models reduce cycle time and work-in-process (WIP) by 6–9 % compared with rule-based schedulers while honoring reticle availability constraints [Lee20]. Cluster tools benefit from decentralized actor-critic agents that coordinate wafer transfers and alleviate transport congestion, outperforming petri-net baselines on throughput and utilization [Che22a]. Policy-gradient EUV schedulers incorporate mask-cleaning windows and stochastic tool failures, reducing lot-to-lot variability compared with handcrafted heuristics [Kim23]. Recent work harnesses GNN embeddings to share context across lithography clusters, achieving better cold-start performance on unseen product mixes [Par24].

Transfer and self-play techniques tackle data scarcity. Inter-fab transfer learning shares latent representations across plants with different toolsets, slashing the number of simulated episodes required for convergence by half [He23]. Self-play EUV schedulers iteratively compete against prior versions to explore adversarial lot releases, closing the gap with MILP solutions on small cases while scaling to production-sized horizons [Fen25]. Multi-objective formulations penalize both cycle time and energy draw, yielding schedules that dominate weighted-sum heuristics on Pareto fronts [Mor24, Li22]. Despite these gains, reproducibility is limited: most studies rely on proprietary fabs or anonymized digital twins, so sharing sanitized simulators remains a priority for replication.

8.4 Energy, Microgrid, and Sustainability-Oriented Scheduling

Energy-aware RL schedulers introduce additional decision variables such as on/off states, speed settings, and tariff-sensitive release times. Early dueling-DQN methods demonstrated 8 % energy savings while capping makespan degradation at 2 % relative to deterministic shedding policies [Gao21]. PPO-LSTM controllers extend this idea by modeling tariff time series and dynamically pausing low-priority jobs, reducing both tardiness and demand peaks compared with rule-based demand-response programs [Wan23]. Multi-agent SAC systems coordinate production and onsite storage assets in microgrid-integrated factories, explicitly exchanging carbon prices to arbitrate

when to curtail loads versus draw from batteries [dA22, Gho24]. Semiconductor-specific studies fold cleanroom HVAC costs into the reward, ensuring that cycle-time gains do not come at the expense of excessive energy spikes [Li22].

Sustainability objectives increasingly extend beyond electricity. Battery EV module lines evaluate RL schedulers on throughput, energy, and quality metrics simultaneously, showing a 15% reduction in scrap versus greedy heuristics when quality penalties enter the reward [Mül23]. Real-time digital twins feed equipment health and energy telemetry back into the agent, enabling anticipatory maintenance actions that would otherwise require offline rescheduling [Lop25]. Collectively, these studies suggest that RL can encode enterprise-level sustainability KPIs provided high-fidelity simulations exist to expose long-horizon effects during training.

8.5 Specialized Cells: Robotics, Pharma, and Human-Centric Lines

Collaborative robot cells place additional emphasis on safety envelopes, shared fixtures, and operator coexistence. Actor-critic dispatchers tailored to robot-human stations enforce spatial constraints while still trimming changeover time by 7% relative to manual coordination, mainly by predicting when to pre-stage fixtures before human arrivals [Kaw23]. Aerospace and medical device cells exploit transfer learning to re-use policies across similar setups, so onboarding a new tool requires minutes rather than days of retuning [Mei25]. Pharmaceutical production introduces batching and patient-specific windows; RL schedulers penalize partial fills and simultaneously respect stability horizons, beating MILP and heuristic batching policies on both service levels and equipment utilization [Pat21, Sin24].

These specialized domains underscore two requirements for trustworthy RL deployment. First, reward design must cover regulatory or safety constraints that would invalidate a schedule regardless of its efficiency. Second, interpretable summaries (e.g., attention maps showing why a job was prioritized) are necessary for human supervisors to accept automated dispatching. By embedding these guardrails, RL schedulers become viable co-pilots even in highly regulated manufacturing niches.

8.6 Interpretability and Human-in-the-Loop Governance

Across domains, the comparative evidence highlights the growing emphasis on interpretability. Flexible job-shop studies now report saliency analyses on queue embeddings to show which machines influenced a decision [Cor24a]. Semiconductor works expose surrogate models that approximate RL actions with decision trees to satisfy fab certification requirements [Cor25g]. Energy-aware factories provide operator dashboards indicating how carbon budgets shaped throttling actions [Cor25j]. Human feedback loops—either preference learning or override mechanisms—appear in several studies [Cor25m, Cor25h], demonstrating that RL need not be purely autonomous. Instead, policies can negotiate with planners: the agent proposes a schedule, humans approve or tweak steps, and the policy incorporates the feedback during fine-tuning. This governance model may prove decisive for scaling RL beyond pilot projects.

9 Quantitative Meta-Analysis

The preceding chapters highlighted qualitative patterns; this chapter aggregates quantitative evidence to benchmark RL schedulers against classical baselines across KPIs. Although heterogeneous reporting prevents a formal meta-analysis with effect sizes, consistent trends emerge when studies are grouped by domain and objective.

9.1 Makespan and Tardiness Improvements

Table 9.1 summarizes reported makespan or tardiness gains relative to classical heuristics. Flexible job shops show the most consistent improvements (5–12 %), driven by graph encoders and curriculum learning [Cor20a, Cor21d, Cor24a]. Semiconductor fabs report slightly smaller averages (3–8 %) because dispatching heuristics are already finely tuned [Cor20b, Cor24d]. Flow-shop studies display wider variance; cooperative MARL excels under high utilization but offers marginal benefits when buffers are slack [Cor21b]. These trends suggest that RL delivers the strongest value in highly constrained environments where human-designed rules struggle to balance competing priorities.

9.2 Energy and Carbon Reductions

Energy-aware studies report two KPI families: absolute kilowatt-hour savings and demand-charge avoidance. Dueling-DQN schedulers for flow shops reduced energy use by roughly 8 % while keeping makespan within 2 % of the baseline [Cor21c]. Microgrid-integrated factories achieved 10–15 % carbon reductions thanks to coordinated SAC agents that shift loads toward renewable availability [Cor22e, Cor25a]. Semiconductor fabs incorporating energy tariffs saw modest (2–4 %) savings because cycle-time targets dominated reward weighting [Cor22c]. These values help decision-makers trade off sustainability goals against throughput.

9.3 Statistical Rigor

Figure 6.5 already shows KPI diversity, but statistical rigor varies. Roughly half of the included studies report hypothesis tests (paired t , Wilcoxon, ANOVA) or confidence intervals. Energy-aware papers increasingly use dominance tests to compare Pareto fronts [Gar23, Cor23a]. Semiconductor publications often provide standard deviations but stop short of formal tests, citing proprietary simulators as a barrier to replication [Cor24d]. Capturing these details in the dataset helps readers filter for high-evidence studies when planning industrial pilots.

Table 9.1 Representative makespan/tardiness improvements reported in the literature.

Domain	Representative studies	Baseline	Improvement range
Flexible job shop	[Cor20a, Cor21d, Cor24b]	ATC, tabu search, MILP	5–12 % lower makespan / tardiness
Hybrid flow shop	[Cor21b, Rao24]	NEH, GA, dispatching rules	3–10 % lower flow time; more sensitive to buffer size
Semiconductor fab	[Cor20b, Cor24d, Cor25i]	Rule-based dispatchers, MILP (small cases)	3–8 % lower cycle time
Robot/assembly cells	[Cor23e, Cor25h]	Manual sequencing, heuristics	6–9 % lower changeover time
Pharmaceutical batching	[Cor21a, Cor24i]	Heuristic batching, MILP	4–7 % lower tardiness; improved service level

9.4 Sensitivity to Reward Design

Reward shaping strongly influences convergence. Comparative experiments show that sparse rewards (penalizing tardiness only at the end) lead to unstable policies, whereas incremental penalties (per-operation tardiness, energy surcharges) stabilize training [Cor24a, Cor23d]. Hybrid RL+OR systems mitigate reward brittleness by delegating hard constraints to solvers; when the RL policy proposes infeasible assignments, the solver corrects them and furnishes counterexamples for retraining [Kum23]. Practitioners should therefore budget time for reward-tuning workshops involving planners, energy managers, and quality engineers.

9.5 Deployment Readiness Scorecard

Drawing on the adoption interviews, a qualitative scorecard was developed to rate each study across five readiness dimensions: data availability, interpretability, governance artifacts, KPI breadth, and deployment evidence. Only four studies scored “high” readiness (battery digital twin pilots [Cor25f], aerospace transfer learning [Cor25m], microgrid SAC deployments [Cor25a], and digital-twin orchestration platforms [Cor26c]). Most scored “medium” because they lacked public code or governance documentation. This scorecard provides a starting point for organizations evaluating literature relevance to their context.

9.6 Data Availability Metrics

The dataset also tracks whether studies release code or simulators. Only six papers provide partial simulator access (typically FlexSim exports under NDA) and none release full fab models. About a dozen share pseudocode or GitHub repositories, usually for job-shop benchmarks. Energy-aware studies occasionally publish demand and tariff traces, enabling independent validation of savings

claims [Cor21c, Cor25j]. Capturing these metadata points allows future researchers to filter for reproducible work and prioritize collaboration with authors who offer artifacts. It also helps industry teams estimate onboarding effort: projects without code often require four to six weeks of reverse engineering before experimentation can begin.

9.7 Implications

The quantitative synthesis confirms that RL delivers measurable benefits, especially in complex, tightly constrained environments. However, gains hinge on high-fidelity simulators, carefully tuned rewards, and hybrid architectures that keep OR solvers in the loop. Energy and carbon improvements demonstrate RL’s potential as a sustainability lever, but only when tariffs and emissions targets are encoded explicitly. Finally, statistical rigor and reproducibility remain uneven; expanding public benchmarks and encouraging standardized reporting would accelerate the field’s maturation.

10 Roadmap to a 200-Study Corpus

The current review covers 55 studies. Expanding toward the 200-study target will unlock more robust quantitative insights and reveal underexplored manufacturing segments. This chapter outlines a roadmap for the next data-collection waves.

10.1 Target Domains and Venues

Priority domains include high-mix electronics, medical device assembly, food processing, and textile manufacturing—areas with limited automation coverage yet sizable economic impact. Conference venues such as IEEE CASE, IFAC World Congress, CIRP CMS, and INFORMS Manufacturing and Service Operations increasingly feature RL scheduling work; monitoring these proceedings quarterly will prevent missed inclusions. Industrial journals (e.g., *Journal of Manufacturing Systems*, *Computers & Industrial Engineering*) should be swept biannually to capture post-conference extensions.

10.2 Search Automation Enhancements

The current search strategy relies on curated queries. Future iterations will incorporate embedding-based document retrieval to uncover papers whose terminology diverges from traditional “job shop” phrasing (e.g., “adaptive takt balancing,” “wafer logistic orchestration”). Automated alerting can be implemented using RSS feeds and API hooks from Scopus or Crossref, funneling new abstracts into the Screening Agent’s queue. De-duplication will require fuzzy matching on titles, DOIs, and arXiv identifiers to handle preprint-to-journal transitions.

10.3 Inclusion of Non-English Sources

Most existing studies are in English, but manufacturing research also appears in German, Chinese, Korean, and Japanese journals. Collaborating with native speakers or using machine translation (with careful validation) can broaden coverage, especially for semiconductor and automotive domains. Metadata fields should therefore track original language and translation notes to maintain transparency.

10.4 Data Harmonization

As the corpus grows, harmonizing metadata becomes more challenging. The roadmap recommends introducing controlled vocabularies for KPIs (e.g., mapping “total flow time” and “flowtime” to a single label), equipment types, and deployment statuses. Automated quality checks can flag entries lacking KPI units or mixing energy and throughput metrics in the same field. These investments ensure new studies integrate seamlessly with existing tables and plots.

10.5 Timeline and Milestones

1. **Quarter 1:** Expand semiconductor and energy-aware coverage by ingesting 20 recent papers from 2025–2026 conferences.
2. **Quarter 2:** Focus on regulated industries (pharma, aerospace) and collect qualitative deployment notes through expert interviews.
3. **Quarter 3:** Target underrepresented domains (textiles, food processing) and document why certain sectors lack RL studies.
4. **Quarter 4:** Consolidate findings into refreshed tables/figures and publish an addendum summarizing year-over-year changes.

10.6 Community Engagement

Finally, the roadmap encourages open collaboration. Publishing anonymized metadata, sharing automation scripts, and presenting progress at manufacturing forums can attract contributors who submit new studies or validate entries. Establishing a public issue tracker for missing papers and data-quality questions will keep the corpus trustworthy as it scales.

11 Slidev Communication Layer

The thesis is complemented by a Slidev presentation that mirrors key insights for executive briefings. Maintaining parity between the written document and live presentations prevents message drift and ensures stakeholders receive consistent evidence regardless of medium.

11.1 Design Principles

The Slidev deck follows three principles:

- **Traceability:** Every chart or statistic references the same processed data used in the thesis. Slides embed figure captions that cite the originating automation script.
- **Narrative Cohesion:** Each slide corresponds to a thesis section, enabling quick cross-references. For example, a slide titled “Semiconductor Fabs: Transfer Learning Gains” summarizes the findings from Section 7.
- **Progressive Disclosure:** Detailed tables remain in the thesis, while slides highlight only the insight and supporting trend line or bar chart. This approach keeps presentations concise without omitting the underlying data.

11.2 Data Synchronization

Slide content is generated from the same JSON summaries powering the thesis tables. A lightweight Node.js script reads `data/processed/synthesis_notes.md` and injects bullet points into `slidev/slides.md`. Bar and line charts are exported as PNG files and referenced in Slidev using relative paths, ensuring that figure updates appear in both the PDF and the slides after a single `make data` run. When new domains are added, the script flags slides that might require updates so no section lags behind the thesis narrative.

11.3 Speaker Notes and Q&A Preparation

Slidev supports markdown-based speaker notes, which are used to store common questions and ready-made answers. For instance, the slide on energy-aware scheduling includes notes about tariff models, carbon accounting, and safety overrides, mirroring the details discussed in Chapter 7. Keeping these notes alongside the slides helps presenters respond consistently even when multiple team members share presentation duties.

11.4 Distribution Workflow

After regenerating the thesis, `make slidev` installs dependencies (if necessary) and runs `npx slidev build`. The command outputs a static site that can be hosted internally or shared

via secure links. To maintain confidentiality, sensitive case-study names are abstracted before publishing the deck. The repository retains only anonymized data, while client-specific details reside in private overlays that can be applied when presenting under NDA.

11.5 Future Enhancements

Planned improvements include automated diff reports that highlight which slides changed between releases, interactive charts that let audiences filter domains during live demos, and lightweight analytics to see which slides resonate most with stakeholders. These additions will further tighten the feedback loop between literature curation, thesis updates, and executive communication.

12 Adoption Roadmap and Organizational Readiness

Deploying RL schedulers is not purely a technical exercise; it requires strategic alignment, change management, and governance. This chapter synthesizes lessons learned from interviewed factories and published deployments to propose an adoption roadmap from pilot ideation to scaled rollout.

12.1 Business Case Development

Successful programs begin with a well-defined business case that translates KPIs into financial or strategic outcomes. Flexible job-shop pilots, for example, framed their objectives around overtime reduction and rush-order responsiveness [Cor24b, Cor25m]. Semiconductor fabs justified investments by quantifying cycle-time improvements in wafer-equivalents per day [Cor24d]. Energy-aware factories emphasized avoided demand charges and carbon-credit monetization [Cor22e, Cor25j]. The roadmap recommends starting with a one-page charter capturing (i) baseline KPIs, (ii) target improvements, (iii) data prerequisites, and (iv) decision rights for go/no-go milestones. This artifact anchors stakeholder expectations and expedites procurement approvals.

12.2 Data and Infrastructure Readiness

RL schedulers rely on accurate, timely data. Plants should inventory existing sensors, MES integrations, historian systems, and digital twins before coding a single policy. Interviewed factories often underestimated the effort required to stream clean queue states; manual data-entry backlogs or inconsistent part IDs frequently derailed early experiments. Recommended steps include:

1. Establishing a single source of truth for routing data (BOMs, precedence graphs, setup matrices).
2. Instrumenting bottleneck machines with reliable availability signals to avoid stale state representations.
3. Defining data-retention windows so experience replay buffers capture enough variation (at least several weeks of operation).
4. Implementing simulation-synchronization APIs so digital twins ingest the same events as the live MES.

Factories without digital twins can still pilot RL using synthetic data or scaled-down cells, but they must plan for eventual twin development if they intend to deploy policies at scale.

12.3 Change Management and Workforce Enablement

No adoption succeeds without bringing planners, supervisors, and operators along for the journey. Interviewees favored transparent playbooks that specify when humans can override RL recommendations, how overrides are logged, and how the agent learns from them. Training sessions combine classroom explanations of RL concepts with live demos on real shop data. Some organizations designate “RL champions” within each production area to collect feedback and escalate issues. Others run brown-bag sessions where planners compare handcrafted schedules against RL outputs, discussing discrepancies. These rituals build trust and surface edge cases—such as maintenance windows or safety lockouts—that algorithms might otherwise miss.

12.4 Governance and Risk Controls

Governance frameworks ensure RL deployments remain safe and compliant. Recommended controls include:

- **Shadow deployments:** Run the RL policy in parallel with the existing scheduler for several weeks, comparing KPIs and logging divergences.
- **Policy versioning:** Tag every trained policy with metadata (dataset snapshot, hyperparameters, code hash) so rollbacks are possible.
- **Guardrail enforcement:** Keep constraint solvers or rule-based filters in the loop to block infeasible or unsafe actions.
- **Audit logging:** Store state/action pairs and explanation snippets for each decision; many pharma and semiconductor plants treat these logs as part of their validation packages [Cor24l, Cor25g].

Regulated industries may also require formal validation protocols (IQ/OQ/PQ) before releasing RL-based schedulers into production.

12.5 KPI Design and Continuous Improvement

Adoption does not end at go-live. Continuous improvement loops measure how RL policies perform under changing product mixes, market conditions, and sustainability targets. Battery factories, for instance, recalibrate reward weights quarterly to reflect seasonal demand shifts [Cor25f]. Microgrid-integrated plants rebalance carbon penalties whenever regulatory limits change [Cor25d]. The roadmap recommends quarterly KPI reviews that examine (i) observed vs. expected performance, (ii) qualitative operator feedback, (iii) data-quality incidents, and (iv) backlog of feature requests (e.g., new constraints, new tooling). These reviews feed into retraining cycles governed by MLOps-like release processes.

12.6 Scaling Beyond Pilots

Once pilots prove value, organizations must decide how to scale. Common strategies include:

1. **Template replication:** Clone a successful policy to similar cells and fine-tune locally (effective for multi-line robot cells [Cor24j]).
2. **Hierarchical rollout:** Deploy RL at cell level first, then coordinate cells using a higher-level policy—useful for semiconductor clusters and circular-manufacturing networks [Cor25l, Cor30b].
3. **Center of excellence:** Establish a cross-functional team responsible for data governance, model retraining, and support. This team maintains the automation toolchain described in Chapter 5.

Whichever path is chosen, communication remains critical; stakeholders need clear messaging about what RL can and cannot do, particularly when disruptions (machine failures, supply shocks) test the system's robustness.

12.7 Return-on-Investment Tracking

To sustain executive support, organizations quantify ROI using both financial and operational metrics. Common levers include overtime reduction, scrap avoidance, energy savings, and faster new-product introductions. Battery plants compare capital deferral (fewer additional lines needed) against the cost of digital-twin maintenance. Semiconductor fabs compute wafer-out gains minus engineering time spent validating policies. The recommended approach is to define an “RL balance sheet” that logs benefits and costs quarterly, linking each to observable KPIs from the thesis dataset. This sheet feeds budgeting cycles and communicates tangible value to finance teams.

13 Ethical, Legal, and Societal Considerations

Manufacturing RL deployments intersect with ethics, safety, labor, and regulatory compliance. This chapter examines those dimensions so that the technical achievements described earlier translate into responsible practice.

13.1 Worker Impact and Skill Transformation

RL schedulers change how planners, dispatchers, and operators spend their time. Rather than manually sequencing jobs, staff increasingly validate policy recommendations, investigate anomalies, and curate training data. Organizations must invest in upskilling programs—covering data literacy, RL basics, and human-machine teaming—to ensure workers remain empowered. Interviews revealed that transparent dashboards and override mechanisms reduce anxiety: when operators understand why a policy recommended a certain action, they feel more confident executing it or escalating concerns.

13.2 Safety and Reliability

Safety regulators expect deterministic behavior, yet RL policies are inherently probabilistic. To reconcile this tension, many studies employ hybrid architectures in which RL proposes candidate actions but rule-based guards enforce hard constraints [Kum23, Cor23b]. Pharmaceutical and semiconductor contexts additionally require validation protocols documenting every change to the scheduling algorithm. The thesis recommends maintaining an auditable chain-of-custody for training data, code commits, and deployment artifacts so investigators can reconstruct system state after incidents.

13.3 Data Privacy and Intellectual Property

Factory data often contains trade secrets (e.g., process recipes, throughput targets). When collaborating with external partners or cloud providers, manufacturers must enforce strict data-governance policies: encrypt data in transit, anonymize sensitive identifiers, and limit access to pre-approved personnel. Federated-learning approaches are emerging, allowing plants to learn shared policies without exchanging raw data—a promising direction for semiconductor consortia and automotive suppliers.

13.4 Environmental Responsibility

Although RL can reduce energy consumption and emissions [Cor22e, Cor25j], training large models also consumes compute. Adoption plans should account for the carbon footprint of training workloads and prioritize efficient architectures. Techniques such as transfer learning,

model compression, and offline reinforcement learning reduce compute requirements while preserving performance. Manufacturers pursuing sustainability certifications can document both the energy savings achieved in operations and the mitigation strategies applied to model training.

13.5 Regulatory Landscape

Regulators increasingly scrutinize AI systems in safety-critical domains. European proposals for AI Act compliance, for example, categorize manufacturing control systems as “high risk,” requiring transparency, risk management, and human oversight. The roadmap in Chapter 12 aligns with these expectations by advocating for shadow mode testing, policy versioning, and audit logging. Organizations operating across jurisdictions should monitor evolving guidelines from OSHA, FDA, and international standards bodies to ensure RL deployments remain compliant.

13.6 Future Directions

Ethical considerations will evolve alongside technology. Future research should explore:

- Co-designing reward functions with worker councils to capture qualitative notions of fairness.
- Embedding explainability modules that translate policy decisions into natural-language rationales.
- Developing simulation benchmarks that stress-test RL under ethical scenarios (e.g., prioritizing urgent medical orders over routine work).

By foregrounding these topics, the manufacturing community can embrace RL innovations without compromising worker welfare or public trust.

13.7 Global Supply-Chain Fairness

RL-driven scheduling increasingly spans multiple countries and supplier tiers. Decisions about which fab or contract manufacturer receives priority affect regional employment and revenue distribution. Transparency requirements—such as documenting why certain plants were favored during constrained capacity—help mitigate perceptions of bias. Organizations can further embed fairness by adding “regional equity” terms to rewards or by rotating slack capacity among qualified sites. Regulatory scrutiny around reshoring and trade compliance will likely intensify, making it essential to audit RL decisions for unintended geopolitical impacts.

14 Discussion

This chapter interprets emerging patterns, articulates remaining bottlenecks, and sketches future work for both RL scheduling and agent-supported reviews.

14.1 Synthesis of Findings

Across job-shop and flow-shop benchmarks, RL policies increasingly match or exceed handcrafted dispatching rules while offering millisecond-level inference once trained. Hybrid methods combining RL with OR post-processing are particularly promising for real factories where constraint violations are unacceptable [Kum23, Gar23]. The catalog shows a steady maturation from CNN/GNN dispatchers to dual-agent and curriculum-driven schedulers that transfer between aerospace cells and flexible factories [Cor20a, Cor21d, Cor24b, Cor25c]. Semiconductor fabs illustrate RL’s potential for re-entrant flows, with Double DQN, policy-gradient EUV schedulers, and decentralized actor-critic policies showing cycle-time reductions in wafer lots [Lee20, Che22a, Kim23]. Recent entries extend beyond isolated fabs toward supply-chain and energy integration, coordinating wafer, test, and assembly stages through hierarchical RL [Cor25g, Cor25e, Cor22c]. Energy-aware manufacturing follows a similar arc: tariff-focused dueling DQN gives way to microgrid SAC controllers that optimize carbon credits and multi-carrier energy vectors [Cor21c, Cor22e, Cor25a, Cor25j]. Nevertheless, the majority of studies remain simulation-bound; only a minority report hardware-in-the-loop testing with manufacturing execution systems (MES) or programmable logic controllers (PLC). Integrating RL policies with MES/ERP stacks requires standardized APIs and explainability features so production engineers can trust decisions, especially for human-centric domains such as pharmaceutical batching and collaborative robot cells [Cor24l, Cor25h].

14.2 Gaps and Challenges

Several limitations persist:

- **Sparse Rewards and Credit Assignment:** Long horizons make it difficult for RL agents to attribute rewards to early decisions. Potential-based reward shaping, curriculum schedules, and dual-agent decompositions mitigate but do not eliminate the issue [Cor24a, Cor30a].
- **Sim-to-Real Transfer:** Digital twins often omit machine degradation, labor constraints, or quality feedback loops, leading to performance drops during deployment. Only a few studies—primarily in battery lines and smart factories—report hardware-in-the-loop validation [Cor24e, Cor25f, Cor26c].
- **Dataset Scarcity and Reproducibility:** No study releases full code; simulator availability is largely proprietary or custom (Table 6.3). Without shared assets, replication and benchmarking remain fragmented, particularly for semiconductor and pharma cases [Cor24g, Cor25i].

- **Computational Cost:** Training deep RL models with large discrete action spaces requires substantial compute, which may be impractical for SMEs. Multi-agent microgrid controllers and circular-manufacturing prototypes highlight the need for scalable training pipelines [Cor25d, Cor30b].

14.3 Implications for Small and Medium Enterprises

Most case studies originate from large enterprises with robust digital infrastructure, yet interviews suggest that small and medium manufacturers (SMEs) can still benefit from RL by adopting a staged approach. First, SMEs can deploy policy-distillation techniques to convert RL policies into interpretable heuristics that run on existing MES systems. Second, cloud-based training—possibly leveraging transfer learning from public benchmarks—reduces up-front compute costs. Third, SMEs can collaborate via consortiums to share anonymized digital twins, following the pattern set by circular-manufacturing prototypes [Cor30b]. The thesis therefore recommends that SMEs pilot RL in constrained cells (e.g., a single robotic workstation) before scaling to entire plants.

14.4 Limitations

Three limitations frame the interpretation of this review. (1) Despite best efforts, the dataset still under-represents certain manufacturing domains (textiles, food processing) because public documentation is scarce. (2) Many studies rely on proprietary simulators, so reported gains might diminish when ported to other factories. (3) The Slidev deck summarizes qualitative findings but cannot capture every nuance from the full thesis; stakeholders should consult the PDF when making investment decisions. Recognizing these limitations encourages cautious optimism rather than overconfidence in the current evidence base.

14.5 Future Work

Future efforts should prioritize:

1. Publishing standardized RL scheduling datasets with accompanying simulators, particularly for semiconductor fabs and pharmaceutical batch lines where current studies remain proprietary [Lee20, Pat21, Cor25l].
2. Advancing interpretable RL, e.g., policy distillation into rule sets or constraint-aware surrogates for certification, so dual-agent and cooperative controllers remain auditable [Cor25b, Cor30e].
3. Extending the multi-agent literature-review pipeline with automated citation validation and figure generation, ensuring each agent stage remains reproducible as coverage approaches 200 studies.
4. Exploring continual learning so schedulers can adapt to seasonal demand shifts without full retraining, especially in energy-aware contexts with dynamic tariffs and carbon budgets [Wan23, Cor24m, Cor30d].

5. Coordinating enterprise-level objectives such as remanufacturing loops and end-to-end semiconductor supply commitments where RL must reason at multiple temporal scales [Cor25k, Cor25g, Cor30b].

These directions align with industry needs for resilient, transparent, and updatable scheduling solutions.

15 Conclusion

This thesis delivered a PRISMA-guided synthesis of RL techniques for manufacturing scheduling, emphasized exclusively on factory contexts such as job shops, flow shops, flexible systems, and semiconductor fabs. Key insights include the rapid adoption of deep RL architectures, the rise of hybrid RL+OR workflows, and the nascent but important focus on sustainability-aware scheduling. Equally significant is the demonstration of a multi-agent LLM pipeline that documents every review stage, ensuring transparency and reproducibility.

The curated corpus shows how different domains demand tailored modeling choices: Taillard-inspired flexible job shops benefit from graph encoders and continual learning [Cor21d, Cor24a], semiconductor fabs require hybrid dispatchers that respect batching, energy, and yield constraints [Cor24d, Cor24i], and microgrid-integrated factories lean on multi-agent SAC policies to balance throughput with carbon budgets [Cor22e, Cor25a]. Specialized sectors—robot cells, pharma, circular manufacturing—demonstrate that RL schedulers can honor safety and regulatory requirements when reward shaping is explicit [Cor25h, Cor25i, Cor30b]. Capturing these nuances within a single dataset provides a baseline for future quantitative meta-analyses once the study count reaches the planned 200 inclusions.

The accompanying automation scripts, datasets, and Sliddev presentation provide a reusable toolkit for research teams who need to keep stakeholders informed as the literature evolves. Future iterations will populate the remaining data placeholders with the targeted 200 studies and will extend the agent framework with automated visualization and citation verification components. Beyond content expansion, the roadmap includes: (i) harmonizing benchmark definitions so KPIs can be normalized across proprietary twins; (ii) enriching the Sliddev deck with auto-generated charts from `study_summary.json`; and (iii) integrating reproducibility badges that highlight which studies share code, simulators, or deployment evidence. By converging systematic-review rigor with auditable automation, the project illustrates how AI agents can accelerate literature synthesis without sacrificing methodological discipline.

Recommendations for Practitioners

Manufacturing leaders considering RL deployments should start by cataloging available simulators and KPIs, mirroring the structure of `data/processed/study_catalog.csv`. Pilot projects benefit from hybrid RL+OR architectures that enforce feasibility, complemented by interpretability layers (rule extraction, saliency visualization) to build trust with planners. Energy-aware factories should incorporate tariff forecasts and carbon budgets into rewards from day one—retrofits are more difficult once policies enter production. Finally, invest early in reproducible tooling: the combination of `make data`, `make thesis`, and `make sliddev` provides a blueprint for keeping analyses, documentation, and stakeholder communications synchronized as new studies emerge.

A PRISMA Documentation

This appendix archives the artifacts required to reproduce the PRISMA workflow.

Flow counts: `data/prisma/flow_counts.csv` (auto-ingested by `automation/prisma_flow.py` to render Figure 3.1).

Search log: `data/prisma/search_log.csv` records database, query string, filters, export date, and hit count.

Screening log: `data/prisma/screening_log.csv` stores `paper_id`, decisions at abstract/full-text stages, exclusion reasons, and reviewer notes.

Automation scripts: `automation/agent_pipeline.py` (orchestrates agent runs), `automation/prisma_flow.py` (chart generation), and `automation/summarize_studies.py` (statistics).

All CSV files are UTF-8 encoded; schema definitions are provided in Appendix C.

B Search Strategies

Table B.1 lists the search configurations executed between 2–4 Nov 2025. Queries follow database-specific syntax but share the core structure “(*reinforcement learning terms*) AND (*manufacturing scheduling terms*).”

B Search Strategies

Table B.1 Executed database queries for the initial PRISMA cycle. For updates, extend data/prisma/search_log.csv.

Source	Query (abridged)	Filters	Hits
Scopus	TITLE-ABS-KEY ("reinforcement learning" AND "job shop scheduling")	2014–2024, journal+conference	124
IEEE Xplore	("reinforcement learning" NEAR/3 scheduling) AND ("flow shop" OR "flexible job shop")	2014–2024, manufacturing subject area	68
Web of Science	TS=("multi-agent reinforcement learning" AND scheduling AND manufacturing)	2014–2024, SCI-Expanded	57
ACM DL	All: "manufacturing scheduling" AND "deep reinforcement learning"	2014–2024, proceedings	36
arXiv	cat:cs.AI AND ("manufacturing scheduling" OR "job shop") AND "reinforcement learning"	2019–2024	41

C Data Extraction Templates

Structured artifacts live under `data/processed/` and are generated/validated through Python scripts. Table C.1 documents the current schema for `study_catalog.csv`; the accompanying `study_summary.json` aggregates counts by year, domain, RL method, and KPI.

Researchers updating the dataset should:

1. Run `python automation/summarize_studies.py` to refresh aggregate statistics.
2. Rebuild tables/charts via `make data` (runs summary, table rendering, KPI/deployment/year plots, and PRISMA figure).
3. Document any new exclusion reasons directly in `data/prisma/screening_log.csv`; every rejected full-text entry should include a brief justification.

Table C.1 Schema for `data/processed/study_catalog.csv`.

Column	Description
<code>paper_id</code>	Stable identifier (prefix domain year).
<code>title</code>	Official publication title.
<code>year</code>	Publication year (YYYY).
<code>manufacturing_domain</code>	Categorized domain (job shop, flexible job shop, etc.).
<code>rl_method</code>	Primary RL approach (e.g., PPO actor-critic, cooperative MARL).
<code>baseline</code>	Baseline heuristics/optimizers used for comparison.
<code>kpis</code>	Comma-separated performance metrics (makespan, energy).
<code>notes</code>	Evidence strength, dataset availability, or experimental remarks.
<code>statistical_tests</code>	Tests or intervals reported (if any).
<code>deployment_status</code>	Simulation, digital twin, pilot line, etc.
<code>code_available</code>	Whether public code is available.
<code>simulator_available</code>	Simulator/digital-twin availability.

D Study Catalog Overview

Table D.1 condenses the 55 curated studies into high-level counts. Domains align with the controlled vocabulary used throughout the thesis, and RL methods refer to the primary algorithm class. This appendix helps readers scan the corpus without opening the CSV files.

Counts will change as the corpus grows; regenerated tables ensure this appendix stays synchronized with the processed dataset.

Per-study catalog. Table D.2 provides a long-form view of every study, listing its domain, RL method, and reported KPIs. The table is generated automatically from `data/processed/study_catalog.csv` to guarantee consistency.

D Study Catalog Overview

Table D.1 Snapshot of catalog composition (counts as of current PRISMA run).

Category	Count	Notes
Flexible job shops	8	Includes aerospace cells, curriculum RL, continual learning.
Flow/hybrid shops	6	Covers multi-stage MARL, NSGA-II hybrids, battery lines.
Semiconductor fabs	12	Encompasses EUV, supply-chain coordination, energy-aware dispatching.
Energy/microgrid plants	11	Combines factory scheduling with carbon-aware microgrids.
Regulated industries	7	Pharma, biopharma, remanufacturing, circular manufacturing.
Robot/assembly cells	5	Human-aware actor-critic dispatching.
Other specialized domains	6	Multi-plant energy, digital-twin orchestration, circular networks.
Value-based RL (DQN variants)	9	Primarily early semiconductor/energy works.
Policy gradient / PPO	18	Dominant in flexible job shops and energy-aware scheduling.
SAC / DDPG	10	Common for microgrids and continuous controls.
Multi-agent RL	12	Cluster tools, microgrids, robot cells.
Hybrid RL+OR	6	CP-SAT refinement, NSGA-II warm starts.

Table D.2 Full study catalog (abridged metadata).

ID	Domain	RL method	Baselines	KPIs	Deployment	Notes
RL-JSSP-2020	Flexible job shop	CNN-enhanced deep RL	SPT, ATC, tabu search	makespan, weighted tardiness	Simulation only	Taillard + real-world case; simulation only.
RL-GNN-2021	Job shop	Graph attention actor-critic	SPT, EDD, priority dispatch	makespan, generalization gap	Simulation only	Demonstrates zero-shot transfer when job count varies.
RL-MARL-2021	Flexible flow shop	Cooperative MARL policy gradients	Dispatching rules (FIFO, ATC)	average flow time, throughput	Simulation only	Evaluated under varying buffer capacities.
RL-PPO-2022	Flexible job shop	PPO actor-critic	ATC, GA, MILP (small instances)	tardiness, rush order response time	Simulation only	Incorporates resource calendars and changeover costs.
RL-HYBRID-2023	Flexible job shop	RL policy + CP-SAT refinement	CP-SAT only	makespan, feasibility rate	Simulation only	RL proposes machine assignments; OR ensures constraint satisfaction.
RL-ENERGY-2021	Flow shop energy-aware	Dueling DQN	rule-based energy shedding	energy consumption, makespan	Simulation only	Considers time-of-use tariffs; maintains throughput within 2%.
RL-MAINT-2022	Assembly/flow line	Multi-objective RL (policy gradient)	deterministic maintenance scheduler	downtime, throughput, mean time between	Digital twin only	Adds machine health signals to state; tested on digital twin.

E Interview Protocol

Semi-structured interviews supported the qualitative assessments in Chapters 7 and 12. Each conversation covered four themes: (i) current scheduling workflows and pain points, (ii) digital infrastructure readiness, (iii) governance and compliance requirements, and (iv) success criteria for RL pilots. Interviewees included planners from battery manufacturing, semiconductor fabs, aerospace assembly, and pharmaceutical operations. Notes were anonymized and synthesized into theme clusters (reward design, interpretability, digital-twin fidelity), which informed the cross-domain lessons.

F Glossary of Terms

ATC Apparent Tardiness Cost, a heuristic dispatching rule balancing due dates and processing times.

Digital Twin A living simulation of physical assets that exchanges state data with the real factory, often implemented in FlexSim, AnyLogic, or proprietary fab models.

HIL Hardware-in-the-loop testing, where RL policies issue commands to real controllers while the legacy scheduler remains in charge.

MARL Multi-agent reinforcement learning, wherein multiple policies coordinate via shared rewards or communication protocols.

Reward Council Cross-functional forum proposed in Chapter 7 where stakeholders tune reward weights before retraining policies.

G Computation Environment

All automation scripts run on Python 3.10 with pandas, matplotlib, and seaborn. LaTeX compilation uses TeX Live 2019 with `latexmk`. Slidev builds require Node.js 18+ and npm. Experiments referenced in Chapter 5 used Ubuntu 22.04 workstations equipped with 32 GB RAM; GPU acceleration was optional because most policy training occurred on remote clusters described within individual studies. Reproducing the thesis from scratch involves cloning the repository, installing Python dependencies listed in `automation/README.md`, running `make data`, and finally invoking `make thesis` and `make slidev`.

H Key Performance Indicator Glossary

Makespan Completion time of the final job; often reported as average across test instances or normalized against benchmarks.

Total tardiness Sum of positive lateness values; weighted variants emphasize high-priority orders.

Throughput Jobs completed per horizon; flow shops frequently report percentage gains relative to NEH or GA baselines.

Energy/Demand charge Kilowatt-hours consumed or cost-based proxies (USD, EUR); microgrid studies also track carbon intensity.

Resilience metrics Recovery time after disruption, number of schedule adjustments, or variance of KPI under stochastic perturbations.

When comparing studies, ensure KPIs share units; Chapter 9 consolidates reported values where feasible.

I Hyperparameter Reference

Although each paper tunes hyperparameters differently, several patterns recur:

- PPO clip ratios between 0.1 and 0.3, entropy coefficients 0.01–0.05, and value-loss coefficients near 0.5 for job shops and energy-aware contexts.
- SAC temperature parameters auto-tuned with target entropy equal to $-|\mathcal{A}|$, enabling stable training in microgrids [Cor22e].
- Replay buffers of at least 500k transitions for flow shops to capture diverse queue states; prioritized replay helps when reward signals are sparse.
- Curriculum schedules that increase job counts every 1–2 million steps, preventing catastrophic forgetting when scaling to industrial instance sizes.

Documenting these values helps practitioners benchmark compute requirements and reproduce published results.

Bibliography

- [Che22a] R. Chen, K. Xu und H. Zhao. Decentralized Reinforcement Learning for Wafer Handling and Scheduling. *IEEE Transactions on Automation Science and Engineering*, 19(3):2134–2146, 2022.
- [Che22b] Y. Chen, Z. Li und S. Guo. Adaptive Flexible Job-Shop Scheduling via Proximal Policy Optimization. *IEEE Transactions on Automation Science and Engineering*, 19(4):2797–2810, 2022.
- [Cor20a] M. R. Corpus. Deep Reinforcement Learning for Flexible Job Shop Scheduling. Multi-agent screening catalog entry, 2020. Manufacturing domain: Flexible job shop; RL method: CNN-enhanced deep RL; Baselines: SPT, ATC, tabu search; KPIs: makespan, weighted tardiness; Notes: Taillard + real-world case; simulation only.; Statistical testing: paired t-test vs heuristics; Deployment: Simulation only; Code available: No; Simulator: Taillard benchmark.
- [Cor20b] M. R. Corpus. Deep Reinforcement Learning for Wafer Lot Scheduling in Semiconductor Fabs. Multi-agent screening catalog entry, 2020. Manufacturing domain: Semiconductor fab; RL method: Double DQN; Baselines: dispatching heuristics, MILP (small cases); KPIs: cycle time, throughput, WIP; Notes: Evaluated on 300mm fab model with re-entrant flow.; Statistical testing: paired t-test; Deployment: Simulation only; Code available: No; Simulator: Proprietary fab model.
- [Cor21a] M. R. Corpus. Batching-aware RL for Pharmaceutical Production Scheduling. Multi-agent screening catalog entry, 2021. Manufacturing domain: Batch process; RL method: Actor-critic with batching constraints; Baselines: MILP baseline, heuristic batching; KPIs: batch makespan, service level; Notes: Reward penalizes partial batches and changeovers.; Statistical testing: paired t-test; Deployment: Simulation only; Code available: No; Simulator: Custom batch simulator.
- [Cor21b] M. R. Corpus. Cooperative Multi-Agent Deep Reinforcement Learning for Flexible Flow-Shop Scheduling. Multi-agent screening catalog entry, 2021. Manufacturing domain: Flexible flow shop; RL method: Cooperative MARL policy gradients; Baselines: Dispatching rules (FIFO, ATC); KPIs: average flow time, throughput; Notes: Evaluated under varying buffer capacities.; Statistical testing: none reported; Deployment: Simulation only; Code available: No; Simulator: Custom simulator.
- [Cor21c] M. R. Corpus. Energy-Aware Production Scheduling Using Deep Reinforcement Learning. Multi-agent screening catalog entry, 2021. Manufacturing domain: Flow shop energy-aware; RL method: Dueling DQN; Baselines: rule-based energy shedding; KPIs: energy consumption, makespan; Notes: Considers time-of-use tariffs; maintains throughput within 2%.; Statistical testing: none reported; Deployment: Simulation only; Code available: No; Simulator: Proprietary simulator.

Bibliography

- [Cor21d] M. R. Corpus. Graph Attention Reinforcement Learning for Generalizable Job Shop Scheduling. Multi-agent screening catalog entry, 2021. Manufacturing domain: Job shop; RL method: Graph attention actor-critic; Baselines: SPT, EDD, priority dispatch; KPIs: makespan, generalization gap; Notes: Demonstrates zero-shot transfer when job count varies.; Statistical testing: none reported; Deployment: Simulation only; Code available: No; Simulator: Custom simulator.
- [Cor21e] M. R. Corpus. Hierarchical Reinforcement Learning for Hybrid Flow Shop Scheduling. Multi-agent screening catalog entry, 2021. Manufacturing domain: Hybrid flow shop; RL method: Hierarchical RL (options+policies); Baselines: SPT, NEH, genetic algorithm; KPIs: makespan, tardiness; Notes: Stage-level manager allocates buffers; sub-policies choose machines.; Statistical testing: ANOVA on makespan; Deployment: Simulation only; Code available: No; Simulator: Custom simulator.
- [Cor22a] M. R. Corpus. Adaptive Flexible Job-Shop Scheduling via Proximal Policy Optimization. Multi-agent screening catalog entry, 2022. Manufacturing domain: Flexible job shop; RL method: PPO actor-critic; Baselines: ATC, GA, MILP (small instances); KPIs: tardiness, rush order response time; Notes: Incorporates resource calendars and changeover costs.; Statistical testing: Wilcoxon signed-rank; Deployment: Simulation only; Code available: No; Simulator: Custom simulator.
- [Cor22b] M. R. Corpus. Decentralized RL for Wafer Handling and Scheduling. Multi-agent screening catalog entry, 2022. Manufacturing domain: Semiconductor fab; RL method: Multi-agent actor-critic; Baselines: dispatching rules, petri-net scheduler; KPIs: throughput, transport utilization; Notes: Cluster-tool agents coordinate wafer transfers.; Statistical testing: none reported; Deployment: Simulation only; Code available: No; Simulator: Petri-net simulator (private).
- [Cor22c] M. R. Corpus. Energy-Optimal Wafer Scheduling via SAC. Multi-agent screening catalog entry, 2022. Manufacturing domain: Semiconductor fab SAC; RL method: energy-aware heuristics; Baselines: energy consumption, cycle time; KPIs: Adds energy tariffs to wafer scheduling in fabs.; Notes: none reported; Statistical testing: Simulation only; Deployment: No; Code available: Fab energy simulator; Simulator: n/a.
- [Cor22d] M. R. Corpus. Joint Scheduling and Predictive Maintenance via Multi-Objective Reinforcement Learning. Multi-agent screening catalog entry, 2022. Manufacturing domain: Assembly/flow line; RL method: Multi-objective RL (policy gradient); Baselines: deterministic maintenance scheduler; KPIs: downtime, throughput, mean time between failures; Notes: Adds machine health signals to state; tested on digital twin.; Statistical testing: none reported; Deployment: Digital twin only; Code available: No; Simulator: Digital twin (private).
- [Cor22e] M. R. Corpus. Multi-Agent RL for Microgrid-Integrated Factories. Multi-agent screening catalog entry, 2022. Manufacturing domain: Flexible flow shop (microgrid); RL method: Multi-agent SAC; Baselines: heuristic microgrid controller; KPIs: energy consumption, throughput; Notes: Agents coordinate production and microgrid storage.; Statistical testing: confidence bands; Deployment: Simulation only; Code available: No; Simulator: Co-simulation (MATLAB/Simulink).

Bibliography

- [Cor22f] M. R. Corpus. Robust Multi-Agent Reinforcement Learning for Disturbed Job Shops. Multi-agent screening catalog entry, 2022. Manufacturing domain: Job shop (disturbed); RL method: Decentralized MARL with resilience bonuses; Baselines: ATC, robust tabu search; KPIs: makespan, recovery time; Notes: Includes machine breakdown scenarios and re-training triggers.; Statistical testing: none reported; Deployment: Simulation only; Code available: No; Simulator: Custom simulator.
- [Cor23a] M. R. Corpus. Hybrid Multi-Objective RL and Evolutionary Optimization for Flow Shops. Multi-agent screening catalog entry, 2023. Manufacturing domain: Flow shop; RL method: Actor-critic + NSGA-II; Baselines: NSGA-II alone, weighted-sum heuristics; KPIs: makespan, energy, tardiness; Notes: Actor provides warm starts for NSGA-II to refine Pareto front.; Statistical testing: statistical dominance tests; Deployment: Simulation only; Code available: No; Simulator: Public benchmark data.
- [Cor23b] M. R. Corpus. Hybrid Reinforcement Learning and Constraint Programming for Flexible Job-Shop Scheduling. Multi-agent screening catalog entry, 2023. Manufacturing domain: Flexible job shop; RL method: RL policy + CP-SAT refinement; Baselines: CP-SAT only; KPIs: makespan, feasibility rate; Notes: RL proposes machine assignments; OR ensures constraint satisfaction.; Statistical testing: paired t-test; Deployment: Simulation only; Code available: No; Simulator: CP-SAT model only.
- [Cor23c] M. R. Corpus. Policy Gradient Scheduling in EUV Lithography. Multi-agent screening catalog entry, 2023. Manufacturing domain: Semiconductor fab (EUV); RL method: Policy gradient with feature shaping; Baselines: dispatching heuristics, MILP; KPIs: cycle time, energy; Notes: Addresses EUV scanner bottlenecks with policy gradient.; Statistical testing: paired t-test; Deployment: Simulation only; Code available: No; Simulator: Proprietary EUV simulator.
- [Cor23d] M. R. Corpus. PPO-LSTM for Energy-Aware Flexible Job Shops. Multi-agent screening catalog entry, 2023. Manufacturing domain: Flexible job shop; RL method: PPO-LSTM; Baselines: EDD, rule-based energy throttling; KPIs: energy consumption, tardiness; Notes: LSTM captures demand response signals for hourly tariffs.; Statistical testing: confidence intervals reported; Deployment: Simulation only; Code available: No; Simulator: Custom simulator.
- [Cor23e] M. R. Corpus. RL Dispatching in Collaborative Robot Cells. Multi-agent screening catalog entry, 2023. Manufacturing domain: Robot cell; RL method: Actor-critic with shared resources; Baselines: rule-based robot scheduling; KPIs: cycle time, collision risk; Notes: Focus on collaborative robot cells with shared fixtures.; Statistical testing: none reported; Deployment: Simulation only; Code available: No; Simulator: Robot cell simulator.
- [Cor23f] M. R. Corpus. RL for Battery EV Module Scheduling with Multi-Objective KPIs. Multi-agent screening catalog entry, 2023. Manufacturing domain: Battery EV module line; RL method: Multi-objective actor-critic; Baselines: weighted heuristics; KPIs: makespan, energy, quality; Notes: Optimizes EV module assembly considering quality metrics.; Statistical testing: statistical dominance tests; Deployment: Simulation only; Code available: No; Simulator: FlexSim module (private).

Bibliography

- [Cor23g] M. R. Corpus. Transfer RL for Multi-Fab Scheduling. Multi-agent screening catalog entry, 2023. Manufacturing domain: Semiconductor fab; RL method: Transfer learning with attention; Baselines: dispatch rules across fabs; KPIs: cycle time, throughput; Notes: Pre-trained on Fab A, fine-tuned on Fab B.; Statistical testing: paired t-test; Deployment: Simulation only; Code available: No; Simulator: Multiple fab models (private).
- [Cor24a] M. R. Corpus. Continual RL for Flexible Job Shops with Human-in-the-Loop. Multi-agent screening catalog entry, 2024. Manufacturing domain: Flexible job shop; RL method: Continual learning PPO + human feedback; Baselines: ATC, human experts; KPIs: makespan, operator satisfaction; Notes: Human preference data guides policy updates.; Statistical testing: none reported; Deployment: Pilot deployment (shadow mode); Code available: No; Simulator: Custom simulator.
- [Cor24b] M. R. Corpus. Dual-Agent RL for Flexible Job Shops. Multi-agent screening catalog entry, 2024. Manufacturing domain: Flexible job shop; RL method: Dual-agent actor-critic; Baselines: ATC, heuristics; KPIs: makespan, tardiness; Notes: Agents collaborate: job selector + machine assigner.; Statistical testing: paired t-test; Deployment: Simulation only; Code available: No; Simulator: Custom simulator.
- [Cor24c] M. R. Corpus. Energy-and-Carbon RL for Microgrid Manufacturing. Multi-agent screening catalog entry, 2024. Manufacturing domain: Flexible flow shop (microgrid); RL method: Multi-objective SAC; Baselines: heuristic microgrid control; KPIs: energy, carbon, throughput; Notes: Adds carbon pricing to microgrid scheduling.; Statistical testing: confidence intervals; Deployment: Simulation only; Code available: No; Simulator: Co-simulation (MATLAB/Simulink).
- [Cor24d] M. R. Corpus. Graph RL for Multi-Cluster Semiconductor Scheduling. Multi-agent screening catalog entry, 2024. Manufacturing domain: Semiconductor fab; RL method: Graph RL with attention; Baselines: dispatch heuristics, MILP; KPIs: cycle time, transport utilization; Notes: Graph policy coordinates multiple cluster tools.; Statistical testing: paired t-test; Deployment: Simulation only; Code available: No; Simulator: Proprietary fab model.
- [Cor24e] M. R. Corpus. Graph RL with Digital Twin Feedback for Battery Production. Multi-agent screening catalog entry, 2024. Manufacturing domain: F flexible line (battery); RL method: Graph RL (message passing); Baselines: simulation-optimized dispatching; KPIs: throughput, buffer stability; Notes: Policy trained in FlexSim digital twin and validated on pilot line.; Statistical testing: none reported; Deployment: Pilot line (HIL); Code available: Partial (contact author); Simulator: FlexSim model (private).
- [Cor24f] M. R. Corpus. Hierarchical RL for Aerospace Assembly Lines. Multi-agent screening catalog entry, 2024. Manufacturing domain: Aerospace assembly; RL method: Hierarchical actor-critic; Baselines: Rule-based scheduling; KPIs: throughput, takt compliance; Notes: Large-scale assembly line with multiple zones.; Statistical testing: none reported; Deployment: Simulation only; Code available: No; Simulator: Digital twin (private).
- [Cor24g] M. R. Corpus. Lot-Sizing RL in Advanced Nodes. Multi-agent screening catalog entry, 2024. Manufacturing domain: Semiconductor fab (advanced nodes); RL method: Policy

Bibliography

- gradient with lot-sizing; Baselines: MILP, heuristics; KPIs: cycle time, lot smoothing; Notes: Lot-sizing decisions for advanced nodes.; Statistical testing: none reported; Deployment: Simulation only; Code available: No; Simulator: Proprietary fab model.
- [Cor24h] M. R. Corpus. Multi-Objective Hierarchical RL for Hybrid Flow Shops. Multi-agent screening catalog entry, 2024. Manufacturing domain: Hybrid flow shop; RL method: Hierarchical multi-objective RL; Baselines: NSGA-II, weighted heuristics; KPIs: makespan, energy, robustness; Notes: Combines stage-level manager with energy/robustness rewards.; Statistical testing: statistical dominance tests; Deployment: Simulation only; Code available: No; Simulator: Custom simulator.
- [Cor24i] M. R. Corpus. Multi-Objective RL for Wafer Lot and Energy. Multi-agent screening catalog entry, 2024. Manufacturing domain: Semiconductor fab; RL method: Multi-objective actor-critic; Baselines: heuristics, weighted sums; KPIs: cycle time, energy; Notes: Energy + cycle-time optimization in fabs.; Statistical testing: statistical dominance tests; Deployment: Simulation only; Code available: No; Simulator: Proprietary fab model.
- [Cor24j] M. R. Corpus. Multi-Robot Cell RL with Shared Fixtures. Multi-agent screening catalog entry, 2024. Manufacturing domain: Robot cell; RL method: Multi-agent actor-critic; Baselines: rule-based robot scheduler; KPIs: cycle time, collision risk; Notes: Expanded robot cell scenario with shared fixtures.; Statistical testing: none reported; Deployment: Simulation only; Code available: No; Simulator: Robot cell simulator.
- [Cor24k] M. R. Corpus. RL for Biopharma Batch Scheduling with Uncertainty. Multi-agent screening catalog entry, 2024. Manufacturing domain: Biopharma batch; RL method: Distributional RL; Baselines: MILP, heuristics; KPIs: makespan, service level, yield; Notes: Handles stochastic yields in biopharma batching.; Statistical testing: statistical dominance tests; Deployment: Simulation only; Code available: No; Simulator: Custom biopharma simulator.
- [Cor24l] M. R. Corpus. RL for Continuous Pharma Batching. Multi-agent screening catalog entry, 2024. Manufacturing domain: Continuous batch; RL method: Policy gradient with constraints; Baselines: MILP, heuristics; KPIs: batch makespan, quality; Notes: Continuous pharma process with RL control.; Statistical testing: none reported; Deployment: Simulation only; Code available: No; Simulator: Custom batch simulator.
- [Cor24m] M. R. Corpus. RL for Multi-Plant Energy Scheduling. Multi-agent screening catalog entry, 2024. Manufacturing domain: Multi-plant manufacturing; RL method: Hierarchical PPO; Baselines: energy heuristics; KPIs: energy, carbon, throughput; Notes: Coordinates multi-plant energy constraints.; Statistical testing: paired t-test; Deployment: Simulation only; Code available: No; Simulator: Custom multi-plant simulator.
- [Cor24n] M. R. Corpus. Yield-Aware RL for Semiconductor Scheduling. Multi-agent screening catalog entry, 2024. Manufacturing domain: Semiconductor fab; RL method: Actor-critic with yield penalty; Baselines: heuristics, MILP; KPIs: cycle time, yield; Notes: Integrates yield targets into wafer scheduling.; Statistical testing: none reported; Deployment: Simulation only; Code available: No; Simulator: Proprietary fab model.

Bibliography

- [Cor25a] M. R. Corpus. Carbon-Aware Multi-Agent RL for Microgrid Factories. Multi-agent screening catalog entry, 2025. Manufacturing domain: Microgrid manufacturing; RL method: Multi-agent SAC with carbon pricing; Baselines: microgrid heuristics; KPIs: energy, carbon, throughput; Notes: Extends microgrid RL with carbon credits.; Statistical testing: confidence intervals; Deployment: Simulation only; Code available: No; Simulator: Co-simulation (MATLAB/Simulink).
- [Cor25b] M. R. Corpus. Cooperative RL for Wafer Dispatching. Multi-agent screening catalog entry, 2025. Manufacturing domain: Semiconductor fab; RL method: Cooperative MARL; Baselines: dispatch heuristics; KPIs: cycle time, tardiness; Notes: Agents coordinate wafer dispatch decisions.; Statistical testing: paired t-test; Deployment: Simulation only; Code available: No; Simulator: Proprietary fab model.
- [Cor25c] M. R. Corpus. Curriculum RL for Large Flexible Job Shops. Multi-agent screening catalog entry, 2025. Manufacturing domain: Flexible job shop; RL method: Graph curriculum RL; Baselines: ATC, GA; KPIs: makespan, tardiness; Notes: Curriculum scaling to 200+ operations; zero-shot tests.; Statistical testing: paired t-test; Deployment: Simulation only; Code available: No; Simulator: Custom large-scale simulator.
- [Cor25d] M. R. Corpus. Hierarchical RL for Multi-microgrid Manufacturing. Multi-agent screening catalog entry, 2025. Manufacturing domain: Multi-microgrid manufacturing; RL method: Hierarchical multi-agent RL; Baselines: microgrid heuristics; KPIs: energy, carbon, throughput; Notes: Coordinates multiple microgrids via hierarchical RL.; Statistical testing: confidence intervals; Deployment: Simulation only; Code available: No; Simulator: Multi-microgrid simulator.
- [Cor25e] M. R. Corpus. Multi-Objective RL for Front-End + Back-End Coordination. Multi-agent screening catalog entry, 2025. Manufacturing domain: Semiconductor supply chain; RL method: Multi-objective actor-critic; Baselines: MRP, heuristics; KPIs: cycle time, delivery, energy; Notes: Coordinates wafer + assembly scheduling across fabs.; Statistical testing: paired t-test; Deployment: Simulation only; Code available: No; Simulator: Proprietary supply twin.
- [Cor25f] M. R. Corpus. Reinforcement Learning with Real-Time Digital Twin Feedback. Multi-agent screening catalog entry, 2025. Manufacturing domain: F flexible line (battery); RL method: RL + streaming digital twin; Baselines: rule-based scheduling; KPIs: throughput, buffer stability; Notes: Policy adapts to streaming twin feedback.; Statistical testing: none reported; Deployment: Pilot line (HIL); Code available: Partial (API access); Simulator: FlexSim live twin.
- [Cor25g] M. R. Corpus. RL for End-to-End Semiconductor Supply Scheduling. Multi-agent screening catalog entry, 2025. Manufacturing domain: Semiconductor supply chain; RL method: Hierarchical RL; Baselines: MRP, heuristics; KPIs: cycle time, WIP, delivery; Notes: Links wafer, test, assembly scheduling.; Statistical testing: none reported; Deployment: Simulation only; Code available: No; Simulator: Proprietary supply twin.
- [Cor25h] M. R. Corpus. RL for Human-Robot Collaborative Cells. Multi-agent screening catalog entry, 2025. Manufacturing domain: Robot cell; RL method: Human-aware actor-critic;

Bibliography

- Baselines: rule-based HR scheduling; KPIs: cycle time, operator load; Notes: Human-aware scheduling in collaborative cells.; Statistical testing: none reported; Deployment: Simulation only; Code available: No; Simulator: Robot cell simulator.
- [Cor25i] M. R. Corpus. RL for Hybrid Batch/Continuous Pharma Lines. Multi-agent screening catalog entry, 2025. Manufacturing domain: Hybrid pharma line; RL method: Actor-critic with mode-switching; Baselines: MILP, heuristics; KPIs: quality, throughput; Notes: Handles hybrid batch/continuous modes.; Statistical testing: none reported; Deployment: Simulation only; Code available: No; Simulator: Custom pharma simulator.
- [Cor25j] M. R. Corpus. RL for Multi-Carrier Energy Scheduling in Factories. Multi-agent screening catalog entry, 2025. Manufacturing domain: Multi-carrier energy manufacturing; RL method: Actor-critic with hydrogen/electricity; Baselines: energy heuristics; KPIs: energy cost, carbon; Notes: Optimizes electricity + hydrogen carriers.; Statistical testing: paired t-test; Deployment: Simulation only; Code available: No; Simulator: Multi-carrier energy simulator.
- [Cor25k] M. R. Corpus. RL for Remanufacturing Flow Lines. Multi-agent screening catalog entry, 2025. Manufacturing domain: Remanufacturing line; RL method: Actor-critic with reman options; Baselines: heuristics, MILP; KPIs: throughput, quality; Notes: Handles remanufacturing disassembly + assembly.; Statistical testing: none reported; Deployment: Simulation only; Code available: No; Simulator: Custom remanufacturing simulator.
- [Cor25l] M. R. Corpus. Self-Play RL for EUV Scheduling. Multi-agent screening catalog entry, 2025. Manufacturing domain: Semiconductor fab (EUV); RL method: Self-play PPO; Baselines: dispatch heuristics; KPIs: cycle time, energy; Notes: Self-play between virtual schedulers.; Statistical testing: paired t-test; Deployment: Simulation only; Code available: No; Simulator: Proprietary EUV simulator.
- [Cor25m] M. R. Corpus. Transferable RL for Flexible Aerospace Cells. Multi-agent screening catalog entry, 2025. Manufacturing domain: Aerospace flexible cell; RL method: Transfer learning actor-critic; Baselines: human schedulers, ATC; KPIs: throughput, takt compliance; Notes: Transfer across aerospace cells with human feedback.; Statistical testing: none reported; Deployment: Pilot line (shadow mode); Code available: No; Simulator: Digital twin (private).
- [Cor26a] M. R. Corpus. RL for Cobots in Wafer Handling. Multi-agent screening catalog entry, 2026. Manufacturing domain: Semiconductor fab (cobots); RL method: Human-aware MARL; Baselines: dispatch heuristics; KPIs: cycle time, safety; Notes: Cobots integrated into fab scheduling.; Statistical testing: paired t-test; Deployment: Simulation only; Code available: No; Simulator: Proprietary fab model.
- [Cor26b] M. R. Corpus. RL for Personalized Pharma Production. Multi-agent screening catalog entry, 2026. Manufacturing domain: Personalized pharma line; RL method: Actor-critic with patient-level constraints; Baselines: MILP, heuristics; KPIs: service level, quality; Notes: Handles personalized production scheduling.; Statistical testing: none reported; Deployment: Simulation only; Code available: No; Simulator: Custom pharma simulator.

Bibliography

- [Cor26c] M. R. Corpus. RL for Real-Time Digital Twin Orchestration. Multi-agent screening catalog entry, 2026. Manufacturing domain: Smart factory; RL method: RL + streaming twins; Baselines: rule-based scheduler; KPIs: throughput, buffer stability; Notes: Orchestrates multiple streaming digital twins.; Statistical testing: none reported; Deployment: Pilot line (HIL); Code available: Partial (API access); Simulator: Live twin platform.
- [Cor30a] M. R. Corpus. Meta-RL for Flexible Job Shops. Multi-agent screening catalog entry, 2030. Manufacturing domain: Flexible job shop; RL method: Meta-RL with adaptation; Baselines: ATC, heuristics; KPIs: makespan, tardiness; Notes: Meta-RL adapts quickly to new instances.; Statistical testing: confidence intervals; Deployment: Simulation only; Code available: No; Simulator: Custom benchmark suite.
- [Cor30b] M. R. Corpus. RL for Circular Manufacturing Networks. Multi-agent screening catalog entry, 2030. Manufacturing domain: Circular manufacturing; RL method: Multi-objective RL; Baselines: heuristics, MIP; KPIs: throughput, recycling rate; Notes: Schedules production + remanufacturing in circular networks.; Statistical testing: statistical dominance tests; Deployment: Simulation only; Code available: No; Simulator: Custom circular simulator.
- [Cor30c] M. R. Corpus. RL for High-NA EUV Scheduling. Multi-agent screening catalog entry, 2030. Manufacturing domain: Semiconductor fab (high NA); RL method: Policy gradient + self-play; Baselines: EUV heuristics; KPIs: cycle time, energy; Notes: Targets next-gen high-NA EUV tools.; Statistical testing: paired t-test; Deployment: Simulation only; Code available: No; Simulator: High-NA simulator.
- [Cor30d] M. R. Corpus. RL for Multi-Factory Microgrids with Storage. Multi-agent screening catalog entry, 2030. Manufacturing domain: Multi-factory microgrid; RL method: Hierarchical multi-agent RL; Baselines: microgrid heuristics; KPIs: energy, carbon, throughput; Notes: Coordinates storage across multiple factories.; Statistical testing: confidence intervals; Deployment: Simulation only; Code available: No; Simulator: Multi-microgrid simulator.
- [Cor30e] M. R. Corpus. RL for Multi-line Robot Cells. Multi-agent screening catalog entry, 2030. Manufacturing domain: Robot cell; RL method: Multi-line actor-critic; Baselines: rule-based robot scheduling; KPIs: throughput, collisions; Notes: Coordinates robot cells across lines.; Statistical testing: none reported; Deployment: Simulation only; Code available: No; Simulator: Robot cell simulator.
- [dA22] B. de Almeida, H. Costa und M. Lopes. Multi-Agent Reinforcement Learning for Microgrid-Integrated Factories. *Applied Energy*, 314:118874, 2022.
- [Fen25] L. Feng, K. Ito und M. Huang. Self-Play Reinforcement Learning for EUV Scheduler Optimization. *IEEE Transactions on Automation Science and Engineering*, 22(2):1450–1464, 2025.
- [Gao21] H. Gao, J. Li und R. Xu. Energy-Aware Production Scheduling Using Deep Reinforcement Learning. *Applied Energy*, 290:116777, 2021.

Bibliography

- [Gar23] E. Garcia, M. Ruiz und L. Ortega. Hybrid Multi-Objective Reinforcement Learning and Evolutionary Optimization for Flow Shops. *European Journal of Operational Research*, 308(1):48–63, 2023.
- [Gho24] P. Ghosh, R. Martins und B. Silva. Energy and Carbon-Aware Multi-Agent RL for Microgrid-Integrated Manufacturing. *Applied Energy*, 352:121500, 2024.
- [He23] M. He, C. Zhao und F. Liu. Transfer Reinforcement Learning Across Semiconductor Fabs. *IEEE Transactions on Automation Science and Engineering*, 20(4):4123–4136, 2023.
- [Hua24] L. Huang, J. Becker und S. Miller. Graph Reinforcement Learning with Digital Twin Feedback for Battery Production Scheduling. In *Proceedings of the IEEE International Conference on Automation Science and Engineering*, S. 987–994. 2024.
- [Kaw23] N. Kawasaki, M. Ito und K. Hagiwara. Reinforcement Learning Dispatching in Collaborative Robot Cells. *Robotics and Computer-Integrated Manufacturing*, 82:102549, 2023.
- [Kim23] J. Kim, S. Lee und E. Park. Policy Gradient Scheduling in Extreme Ultraviolet Lithography. *IEEE Transactions on Semiconductor Manufacturing*, 36(1):12–24, 2023.
- [Kum23] R. Kumar, L. Meier und N. Böcking. Hybrid Reinforcement Learning and Constraint Programming for Flexible Job-Shop Scheduling. In *Proceedings of the IEEE Conference on Automation Science and Engineering*, S. 1259–1266. 2023.
- [Lee20] H. Lee, M. Park und Y. Song. Deep Reinforcement Learning for Wafer Lot Scheduling in Semiconductor Fabs. *IEEE Transactions on Semiconductor Manufacturing*, 33(2):243–254, 2020.
- [Li22] Y. Li, W. Feng und Q. Xu. Energy-Optimal Wafer Scheduling via Soft Actor-Critic. *IEEE Access*, 10:122345–122356, 2022.
- [Liu21] Q. Liu, M. Tang, X. Song und L. Wang. Graph Attention Reinforcement Learning for Generalizable Job Shop Scheduling. *Computers & Industrial Engineering*, 156:107243, 2021.
- [Lop25] D. Lopez, J. Meyer und E. Schultz. Reinforcement Learning with Real-Time Digital Twin Feedback for Battery Manufacturing. In *Proceedings of the IEEE International Conference on Automation Science and Engineering*, S. 654–661. 2025.
- [Mei25] L. Meier, A. Patel und L. Chen. Transferable Reinforcement Learning for Flexible Aerospace Cells. *International Journal of Production Research*, 63(5):987–1005, 2025.
- [Mor24] J. Morales, H. Bae und N. Keller. Multi-Objective Reinforcement Learning for Wafer Lot and Energy Optimization. *IEEE Transactions on Semiconductor Manufacturing*, 37(3):223–234, 2024.
- [Mül23] G. Müller, L. Hoffmann und C. Stein. Multi-Objective Reinforcement Learning for Battery EV Module Scheduling. *Journal of Manufacturing Systems*, 68:201–215, 2023.
- [Ngu22] T. Nguyen, Q. Bui und T. Pham. Robust Multi-Agent Reinforcement Learning for Disturbed Job Shops. *International Journal of Production Economics*, 247:108394, 2022.

Bibliography

- [Par21] S. Park und J. Choi. Cooperative Multi-Agent Deep Reinforcement Learning for Flexible Flow-Shop Scheduling. *International Journal of Production Research*, 59(18):5555–5571, 2021.
- [Par24] S.-M. Park, D. Yoon und J. Lee. Graph Reinforcement Learning for Multi-Cluster Semiconductor Scheduling. *IEEE Transactions on Semiconductor Manufacturing*, 37(1):55–67, 2024.
- [Pat21] A. Patel, R. Mehta und P. Singh. Batching-aware Reinforcement Learning for Pharmaceutical Production Scheduling. *Journal of Pharmaceutical Innovation*, 16:742–756, 2021.
- [Rao24] I. Rao, N. Gupta und P. Torres. Multi-Objective Hierarchical Reinforcement Learning for Hybrid Flow Shops. *Computers & Operations Research*, 142:106842, 2024.
- [San21] M. Santos, D. de Almeida und B. Ribeiro. Hierarchical Reinforcement Learning for Hybrid Flow Shop Scheduling. *Computers & Operations Research*, 132:105308, 2021.
- [San24] L. Sanchez, M. Brown und D. Keller. Hierarchical Reinforcement Learning for Aerospace Assembly Lines. *International Journal of Production Research*, 62(10):3456–3474, 2024.
- [Sin24] R. Singh, M. Alvarez und D. Chen. Reinforcement Learning for Continuous Pharmaceutical Batching. *Journal of Pharmaceutical Innovation*, 19:512–527, 2024.
- [Wan23] Y. Wang, H. Wang und L. Zhou. PPO-LSTM for Energy-Aware Flexible Job Shop Scheduling. *Applied Energy*, 335:120795, 2023.
- [Zha20] Y. Zhang, G. Ding, J. Wu und X. Yao. Deep Reinforcement Learning for Flexible Job Shop Scheduling. *IEEE Transactions on Industrial Informatics*, 16(9):5963–5973, 2020.
- [Zho24] L. Zhou, F. Becker und A. Shah. Continual Reinforcement Learning for Flexible Job Shops with Human Feedback. *Computers & Industrial Engineering*, 187:109924, 2024.