# Choice-Based Learning in JAX

**Shangyin Tan**[*][†]   **Dan Zheng**[*]   **Gordon Plotkin**   **Ningning Xie**
Google DeepMind
{shangyin,danielzheng,plotkin,ningningxie}@google.com

## Abstract

Choice-based learning is a programming paradigm for expressing learning systems in terms of choices and losses. We explore a practical and modular implementation of choice-based learning in JAX by combining two techniques in a novel way: algebraic effects and loss effects. We describe the design and implementation of our library, explore its usefulness for real-world applications such as hyperparameter tuning and deep reinforcement learning, and compare it with existing approaches.

## 1   Introduction

Machine learning has achieved many successes in recent decades, but modular and flexible machine learning programming remains a challenge.

*Choice-based learning* [1, 3, 4] is a recently introduced programming paradigm that can address this challenge. In choice-based learning, programs make *choices* and learn from *feedback* in the form of losses and rewards. This paradigm prioritizes modularity and can be used to implement widely used decision-making techniques, including reinforcement learning and gradient descent.

In this paper, we present our progress in implementing CHOIX, a library for choice-based learning in Python. We describe the design and implementation of CHOIX and use it to explore the usefulness of choice-based learning as a modular programming paradigm for machine learning applications.

## 2   Background

In this section, we explain what choice-based learning is and describe how it can be implemented via two programming paradigms: algebraic effects and handlers, and losses and choice-loss continuations.

**Choice-based learning.**   Choice-based learning systems make choices based on losses (equivalently one could work with rewards). This is done via two operations: `choose` and `loss`. `choose` selects a choice from a choice input space, recrds a numerical loss value, representing the feedback resulting from the choices made. Choice strategies make choices intended to minimize accumulated loss, so that, over time systems learn to make better choices.

**Algebraic effects and effect handlers.**   Algebraic effects are a well-explored technique for structured control flow abstraction in programming languages [6, 7]. Programs use abstract effect operations, and effect handlers provide meaning for such effect operations.

Algebraic effects provide a natural way to model choice-based learning programs: `choose` operations can be represented as effect operations and choice strategies as effect handlers. This allows

---

[*]Equal contribution.
[†]Work done during internship at Google DeepMind.

programmers to describe learning algorithms in terms of choices and losses, while leaving choice strategies up to separately-defined optimization algorithms.

Parameterized handlers [7] extend effect handlers by allowing state to be shared between effect operations. This is useful for implementing choice strategies that need to access and update state in a modular way.

**Losses and choice-loss continuations.** Prior work [1] establishes semantics for choice-based learning languages with `loss` operations. In the semantics, `loss` is a built-in effect operation that can be used freely within programs, including inside handlers. Programs containing `loss` operations are transformed to return an accumulated loss value in addition to their original results.

Additionally, effect handlers can access loss values via a *choice-loss continuation* (or just "loss continuation", `lk`). Loss continuations compute the accumulated loss from using a particular choice; effect handlers can make optimization choices using loss continuations by optionally transforming them (e.g., via automatic differentiation) and calling them with different arguments before resuming the program with the optimal choice.

## 3   CHOIX: a library for choice-based learning

In this section, we introduce CHOIX[1], a library for choice-based learning in Python using JAX [2]. CHOIX supports choice-based learning via the following features: user-defined effect operations and effect handlers (as described in Section 2), `loss` operations, and loss continuations.

In CHOIX, effect operations are declared as Python functions decorated with `@effect`: these functions define an abstract interface and need no implementation. Effect handler functions provide concrete implementations for effect operations. CHOIX provides a `handle` function for effect-handling: `handle` takes a list of handlers and an "effectful" function containing effect operations like `choose` and `loss`, and additionally returns an effect-handled version that returns an accumulated loss. With `handle`, users can write programs using rich features for choice-based learning, then convert them to standard JAX functions to be used normally with other JAX code, including transformations and compilation.

### 3.1   Examples

In this section, we walk through machine learning programs written using CHOIX to highlight the benefits of programming in the choice-based learning paradigm.

#### 3.1.1   Linear regression

Linear regression provides a simple application of choice-based learning: given training data $X, Y$ and initial parameters $\theta$, choose new parameters $\theta'$ that minimize a loss function via gradient descent. Figure 1 shows examples of linear regression written in standard JAX and in CHOIX.

With CHOIX, gradient descent is modeled as a `choose` effect that updates parameters: $\theta' = \text{choose}(\theta)$. A loss value $l$ is computed from the updated parameters by comparing predictions with target data, and `loss`($l$) is called to provide feedback. A `choose_grad` handler function implements parameter update: it first gets gradients $g$ by differentiating the loss continuation $lk$ with respect to $\theta$, then resumes the program with updated parameters $\theta' = \theta - \eta g$, where $\eta$ is a fixed learning rate.

Note that we are not presenting the benefit of CHOIX in the linear regression example. Instead, the goal is to demonstrate CHOIX can be used to to create choice-based learning programs. As choice-based learning examples grow more complex (Appendix C, Section 3.1.2, Section 3.1.3), we will highlight how CHOIX enables modular programming, expressing complex choice-based learning applications in a structured manner. For example, Appendix C shows how to extend linear regression with hyperparameter tuning in CHOIX using nested effect operations and `handle`.

---

[1]"Choix" means "choice" in French.

```
1  def loss_fn(params, x, y):                          1   @effect
2      w, b = params                                    2   def choose[T](x: T) -> T:
3      prediction = w * x + b                           3     """Choice effect operation."""
4      return jnp.mean(prediction - y ** 2)             4
5                                                       5   def linear_regression(params, x, y):
6  def update(params, grads, lr):                       6       """Performs one step of linear regression."""
7      return jax.tree_map(                             7       w_new, b_new = choose(params)
8          lambda p, g: p - lr * g, params, grads       8       prediction = w_new * x + b_new
9      )                                                9       loss(jnp.mean((prediction - y) ** 2))
10                                                      10      return (w_new, b_new)
11  def linear_regression(lr, params, x, y):            11
12      """Performs one step of linear regression."""   12  def update(params, grads, lr):
13      loss_value = jax.grad(loss_fn)(params, x, y)    13      return jax.tree_map(
14      grads = grad_fn(params, x, y)                   14          lambda p, g: p - lr * g, params, grads
15      new_params = update(params, grads, lr)          15      )
16      return new_params                               16
17                                                      17  def choose_grad_handler(lr, params, k, lk):
18  # Usage.                                            18      """Gradient-descent handler."""
19  x, y, w, b = ...                                    19      grads = jax.grad(lk, argnums=(0))(params)
20  new_params = linear_regression(                     20      new_params = update(params, grads, lr)
21      lr=0.001, (w, b), x, y                          21      return k(new_params)
22  )                                                   22
                                                        23  # Usage.
                                                        24  linear_regression = handle(
                                                        25      Handler(parameterized=True, choose=choose_grad_handler)
          losses: loss-related operations               26  )(linear_regression)(lr=0.001)
          effects: effect operations                    27
          handlers: effect handlers                     28  x, y, w, b = ...
                                                        29  loss_value, new_params = linear_regression((w, b), x, y)
```

Figure 1: Linear regression via gradient descent in standard JAX (left) versus choice-based learning in CHOIX (right). In `choose_grad_handler`, the continuation `k` resumes the program from where the effect operation is called (line 18).

### 3.1.2  Reinforcement learning

Reinforcement learning (RL) is a particularly good fit for choice-based learning, as effect operations provide a modular interface for RL algorithms, and parameterized effect handlers facilitate seamless state sharing (e.g. agent parameters like Q-tables or neural network weights) between operations.

In Appendix D, we show a generic function for RL in CHOIX called `run_episode`. `run_episode` uses multiple effect operations: agent operations `predict` and `feedback` and an environment operation `observe`. `predict` selects an action based on the current state and agent parameters; `observe` applies an action to the state to get a new state and reward; finally, `feedback` updates agent parameters based on reward and accumulated loss. By implementing handler functions for these operations, we can get a specialized version of `run_episode` for a particular agent algorithm and environment.

Handler functions in Appendix D implement Q-learning [11], a classic RL algorithm that uses a Q-table mapping state-action pairs to expected rewards. The `predict` handler selects an action based on the Q-table and state, and the `feedback` handler updates the Q-table based on the action and reward. The `observe` handler here implements environment update for the cliff-walking task [10].

### 3.1.3  Deep reinforcement learning

Deep reinforcement learning is a particularly compelling variant of reinforcement learning where the RL agent is a deep neural network and the training process aims to update parameters to enhance predictions. Implementing deep reinforcement learning with CHOIX is as simple as reusing the generic RL algorithm from Section 3.1.2 with alternative handlers.

In Appendix E, we present an implementation of proximal policy optimization (PPO) [8], a well-known algorithm in deep reinforcement learning. Compared with the Q-learning example above, the main changes are in the `feedback_handler`. The PPO algorithm updates agent network weights using three losses: actor loss, value loss, and entropy. These loss values are naturally accumulated via `loss` without the need of writing an explicit pure loss function.

For parameter update, we use gradient descent with the Adam optimizer from Optax[2]. This is done by adapting the `choose` effect from the linear regression example and modifying the gradient descent handler to use Optax. Notably, `choose_optax_handler` is generic like `choose_grad_handler`: both can be used in other programs for gradient optimization.

---

[2] https://github.com/google-deepmind/optax

Additionally, minor adjustments to the `predict` and `observe` handlers are required, reflecting changes in the agent parameter and environment. In `ppo_predict_handler`, the action is selected by sampling from a categorical distribution based on agent network logits. In order to train on Atari games, we change the `observe` handler to `atari_observe_handler` to correctly update the Atari environment.

### 3.1.4 Summary

The examples above demonstrate how Choix enables writing modular choice-based learning programs. The separation between abstract effect operations and concrete handler implementations makes the process of developing new choice-based learning programs easier: new programs define their algorithms as abstract effect operations where effect handlers can be shared and composed between programs. For example, most choice-based learning programs that use Optax can reuse `choose_optax_handler` without changing their own application code.

### 3.2 Implementation

JAX[2] is an open-source system for machine learning research. It offers the familiar API of NumPy with the speed of hardware accelerators via the XLA compiler. JAX's design centers around composable function transformations such as automatic differentiation, compilation, and parallelization. This makes it a solid basis for our work.

CHOIX is implemented as a function transformation on JAXPRs[3], the core language underlying JAX. The transformation is done in multiple stages, producing an *effect-handled* JAXPR at the end. Appendix A explains the transformation stages and characteristics of transformed JAXPRs.

## 4 Related work

*Choice-based learning.* Choice-based learning as a programming paradigm has gained attention in recent years[1, 4]. Frameworks like SmartChoices[3] and PyGlove[5] provide abstractions for writing programs in terms of decisions and rewards, but typically provide a fixed search interface and limited search algorithms. CHOIX supports flexible search thanks to effect handlers: users can define custom search operations, implement or import search algorithms from external frameworks like Vizier[9], and compose them together to create modular choice-based programs.

*Effect handling in JAX.* To our knowledge, our work is the first to propose a library for algebraic effects in JAX. Our work was inspired by a simple proof-of-concept for effect handling in JAX[4], but goes beyond it to implement a fully-fledged effects system. Oryx[5] is a library for probabilistic programming built on JAX. Oryx provides transformations for writing stateful computations via collecting and injecting tagged values, which can be seen as a limited effect system particularly useful for writing deep learning and probabilistic programs. In CHOIX, state management is done in a similar way via parameterized handlers: we provide a short explanation in Appendix B.

## 5 Conclusion

Choice-based learning is a programming paradigm in which systems make feedback-guided choices. We introduce CHOIX, a library for choice-based learning in JAX. CHOIX uses (1) algebraic effects and handlers as a modular interface for writing choice-based learning programs, with (2) the selection monad to enable loss-based choice optimization.

We explore several machine learning applications of choice-based learning in CHOIX. We believe it is a promising approach for developing learning programs, including systems applications that replace hardcoded heuristics with learned strategies. We look forward to exploring more applications of this programming model and welcome feedback and ideas.

---

[3] `https://jax.readthedocs.io/en/latest/jaxpr.html`

[4] `https://colab.sandbox.google.com/drive/1HGs59anVC2AOsmt7C4v8yD6v8gZSJGm6`, linked from `http://blog.sigfpe.com`.

[5] `https://github.com/jax-ml/oryx`

# References

[1] Martin Abadi and Gordon Plotkin. "Smart choices and the selection monad". In: *Logical Methods in Computer Science* 19 (2023).

[2] James Bradbury et al. *JAX: composable transformations of Python+NumPy programs*. Version 0.4.17. 2018. URL: http://github.com/google/jax.

[3] Daniel Golovin et al. *SmartChoices: Augmenting Software with Learned Implementations*. 2023. arXiv: 2304.13033 [cs.SE].

[4] Ugo Dal Lago, Francesco Gavazzo, and Alexis Ghyselen. *On Reinforcement Learning, Effect Handlers, and the State Monad*. 2022. arXiv: 2203.15426 [cs.PL].

[5] Daiyi Peng et al. *PyGlove: Symbolic Programming for Automated Machine Learning*. 2021. arXiv: 2101.08809 [cs.LG].

[6] Gordon D. Plotkin and John Power. "Adequacy for Algebraic Effects". In: *Foundations of Software Science and Computation Structures, 4th International Conference, FOSSACS 2001 Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2001 Genova, Italy, April 2-6, 2001, Proceedings*. Ed. by Furio Honsell and Marino Miculan. Vol. 2030. Lecture Notes in Computer Science. Springer, 2001, pp. 1–24. DOI: 10.1007/3-540-45315-6\_1. URL: https://doi.org/10.1007/3-540-45315-6%5C_1.

[7] Gordon D. Plotkin and Matija Pretnar. "Handlers of Algebraic Effects". In: *Programming Languages and Systems, 18th European Symposium on Programming, ESOP 2009, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2009, York, UK, March 22-29, 2009. Proceedings*. Ed. by Giuseppe Castagna. Vol. 5502. Lecture Notes in Computer Science. Springer, 2009, pp. 80–94. DOI: 10.1007/978-3-642-00590-9\_7. URL: https://doi.org/10.1007/978-3-642-00590-9%5C_7.

[8] John Schulman et al. *Proximal Policy Optimization Algorithms*. 2017. arXiv: 1707.06347 [cs.LG].

[9] Xingyou Song et al. *Open Source Vizier: Distributed Infrastructure and API for Reliable and Flexible Blackbox Optimization*. 2023. arXiv: 2207.13676 [cs.LG].

[10] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. Cambridge, MA, USA: A Bradford Book, 2018. ISBN: 0262039249.

[11] Christopher J. C. H. Watkins and Peter Dayan. "Technical Note Q-Learning". In: *Mach. Learn.* 8 (1992), pp. 279–292. DOI: 10.1007/BF00992698. URL: https://doi.org/10.1007/BF00992698.

## Appendix

### A CHOIX transformation details

CHOIX is implemented as a function transformation on JAXPRs[6], the core language underlying JAX. The transformation is done in multiple stages:

- **Source program.** Source Python programs using CHOIX are lowered to JAXPRs that contain multiple constructs not present in standard JAXPRs: effect operations, effect handlers, and `loss` operations. All of these must be translated away before CHOIX programs can be executed. (Figure 2b.)

- **Delimited, concrete handlers.** Effect handlers now have explicitly delimited scopes. This is needed to create properly-scoped continuations during effect handling. Handler functions are transformed from Python functions to JAXPRs via tracing. (Figure 2c.)

- **Loss transformation.** Performs a selection monad [1] transformation: `loss` operations and choice-loss continuations are translated away. The translated programs compute and return accumulated loss values. (Figure 3a.)

- **Effect handling.** Replaces effect operations with their handler implementations, reifying delimited program continuations using handler scopes. In this manner, Effect operations and effect handlers are translated away. (Figure 3b.)

Finally, we obtain CHOIX-transformed JAXPRs as in Figure 3c. These JAXPRs do not contain effect operations, effect handlers, or `loss` operations: they work naturally with JAX transformations (e.g. `jax.grad` and `jax.jit`), allowing them to interoperate with other JAX code. In this way, CHOIX lets users enjoy the benefits of writing choice-based learning programs using algebraic effects without sacrificing performance.

Optimized runtime performance is not a current focus of CHOIX: since CHOIX produces JAXPRs similar to standard JAX code, we expect similar performance during runtime. In practice, we noticed a slowdown of approximately 0-10% for CHOIX on the example programs included in this paper.

This section walks through the CHOIX program transformation for a toy program, highlighting how it works, stage-by-stage.

The toy program is intentionally simple to save space, since transformation outputs are long. Later appendix sections show more sophisticated programs using CHOIX.

---

[6]`https://jax.readthedocs.io/en/latest/jaxpr.html`

```
1   @effect
2   def increment(x: int) -> int:
3       pass
4
5   def increment_handler(x, k, lk):
6       return k(x + 1)
7
8   @handle(Handler(increment=increment_handler))
9   def main(x):
10      y = increment(x)
11      loss(y ** 2)
12      return y
13
14  print(main(10))
```

(a) Source program.

```
1   # Source Jaxpr.
2   { lambda ; a:i32[]. let
3       _:bool[] = handler[
4         handler_impl={'increment': <function increment_handler>}
5         in_tree=PyTreeDef(())
6         name=increment
7         parameterized=False
8       ]
9       b:i32[] = increment[in_tree=PyTreeDef(*) is_effect_op=True name=increment] a
10      c:i32[] = integer_pow[y=2] b
11      d:f32[] = convert_element_type c
12      _:f32[] = loss d
13      _:bool[] = handler_end_marker[name=increment]
14    in (b,) }
```

(b) Source program as a JAXPR.

```
1   # Delimited concrete handler Jaxpr.
2   { lambda ; a:i32[]. let
3       b:i32[] = delimited_handler[
4         args=[]
5         body_jaxpr={ lambda ; c:i32[]. let
6             d:i32[] = increment[in_tree=PyTreeDef(*) is_effect_op=True name=increment] c
7             e:i32[] = integer_pow[y=2] d
8             f:f32[] = convert_element_type[new_dtype=float32 weak_type=False] e
9             g:f32[] = loss f
10          in (d,) }
11        handler_impl=[
12          increment={ lambda ; h:i32[]. let
13              i:i32[] = add h 1
14              j:i32[] = call_k_with_loss[name=increment] i
15            in (j,) }
16        ]
17        in_tree=PyTreeDef(())
18        name=increment
19        parameterized=False
20      ] a
21    in (b,) }
```

(c) Delimited concrete handler JAXPR.

Figure 2: Initial CHOIX program transformation stages. A Python source program using CHOIX is lowered to a JAXPR. Then, effect handler scopes are made explicit and handler functions are converted from Python to JAXPRs.

7

```
1    # Loss-translated Jaxpr.
2    { lambda ; a:i32[]. let
3        b_loss:f32[] c:i32[] = delimited_handler[
4          args=[]
5          body_jaxpr={ lambda ; d:i32[]. let
6              e:i32[] = increment[in_tree=PyTreeDef(*) is_effect_op=True name=increment] d
7              f:i32[] = integer_pow[y=2] e
8              g:f32[] = convert_element_type[new_dtype=float32 weak_type=False] f
9            in (g, e) }
10         handler_impl=[
11           increment={ lambda ; h:i32[]. let
12               i:i32[] = add h 1
13               j_loss:f32[] k:i32[] = call_k_with_loss[name=increment] i
14             in (j_loss, k) }
15         ]
16         in_tree=PyTreeDef(())
17         name=increment
18         parameterized=False
19       ] a
20     in (b_loss, c) }
```

(a) Loss-translated JAXPR. Effect operations and handlers now return loss values, accumulated from `loss` calls.

```
1    # Effect-handling Jaxpr.
2    { lambda ; a:i32[]. let
3        b_loss:f32[] c:i32[] = closed_call[
4          call_jaxpr={ lambda ; d:i32[]. let
5              h':f32[] i':i32[] = closed_call[
6                call_jaxpr={ lambda ; j:i32[]. let
7                    k:i32[] = add j 1
8                    l_loss:f32[] m:i32[] = closed_call[
9                      call_jaxpr={ lambda ; n:i32[]. let
10                         f:i32[] = integer_pow[y=2] n
11                         g:f32[] = convert_element_type[
12                           new_dtype=float32
13                           weak_type=False
14                         ] f
15                       in (g, n) }
16                   ] k
17                 in (l_loss, m) }
18             ] a
19           in (h', i') }
20       ] a
21     in (b_loss, c) }
```

(b) Effect-handled JAXPR. All special constructs from CHOIX (including effect operations and handlers) have been translated away.

```
1    # Inlined Jaxpr.
2    { lambda ; a:i32[]. let
3        k:i32[] = add a 1
4        f:i32[] = integer_pow[y=2] k
5        g:f32[] = convert_element_type[new_dtype=float32 weak_type=False] f
6      in (g, k) }
```

(c) Inlined, canonicalized JAXPR. This is a standard JAXPR, ready for execution or transformation.

Figure 3: Final CHOIX program transformation stages, performing the loss translation and effect handling. In the final stage, all choice-based learning constructs are translated away, producing a lowered JAXPR that matches JAXPRs from direct-style JAX code.

## B  Stateful computation

```
1   State = int
2
3   @effect
4   def get() -> State:
5       """Gets the state value."""
6
7   @effect
8   def set(s: State) -> None:
9       """Updates the state value."""
10
11  def program() -> State:
12      x1 = get()
13      set(x * 2)
14      x2 = get()
15      set(x2 + 5)
16      return x1 + x2 + get()
17
18  def get_handler(s: State, x: None, k, lk):
19      jax.debug.print('get: s = {}', s)
20      return k(s, s)
21
22  def set_handler(s: State, x: State, k, lk):
23      jax.debug.print('set: s = {}, x = {}', s, x)
24      return k(x, None)
25
26  x = 42
27  main = handle(
28      Handler(parameterized=True, get=get_handler, set=set_handler)
29  )(program)(x)
30
31  result = main()
32  # get: s = 42
33  # set: s = 42, x = 84
34  # get: s = 84
35  # set: s = 84, x = 89
36  # get: s = 89
37  print(result)
38  # (loss=0.0, value=Array(215, dtype=int32))
39
40  print(jax.make_jaxpr(main)())
41  # { lambda ; a:i32[]. let
42  #     b:i32[] = mul a 2
43  #     c:i32[] = add b 5
44  #     d:i32[] = add a b
45  #     e:i32[] = add d c
46  #   in (0.0, e) }
```

Figure 4: A simple stateful program in CHOIX, using `get` and `set` operations. This is noteworthy since state management is a primary concern of deep learning libraries built on JAX, as JAX's purely functional programming model does not provide native support for state management[7]. As shown, the final transformed JAXPR is as efficient as direct-style JAX code.

---

[7]https://jax.readthedocs.io/en/latest/jax-101/07-state.html#taking-it-further

## C Hyperparameter optimization

In machine learning, *hyperparameters* are variables (like learning rate) that are not directly related to the training data but govern the training process itself. Hyperparameter values can greatly affect the training process, including how quickly training converges and whether training succeeds at all. *Hyperparameter optimization* (or tuning) is the process of searching for optimal hyperparameters.

Standard machine learning approaches tend to view hyperparameter tuning as a separate and orthogonal process to model training, often using external blackbox optimization frameworks like Vizier [9] that involve nontrivial configuration. This leads to a barrier between training and tuning, making it difficult to extend programs for training to also do hyperparameter search without rewriting code.

With CHOIX, hyperparameter tuning is naturally represented using *nested effects and handlers*, as shown in Figure 5. We define an effect `choose_single` for selecting a hyperparameter configuration from a search space. This lets us write a `hyperparameter_tuning` generic function that takes a search space and applies `choose_single` to select hyperparameters for an argument training function `f`.

Hyperparameter optimization approaches can then also be implemented as handlers. Here, the `choose_enumerate_handler` function defines a simple grid search that exhaustively explores a search space and picks the option minimizing loss. Other handlers could implement more advanced features such as random exploration or early-stopping. Existing libraries for hyperparameter search like Vizier [9] could also be used to implement handlers.

```python
@effect
def choose_single[T](x: Sequence[T]) -> T:
  """Choose-one-from-many effect operation."""

def hyperparameter_tuning[T](hparam_options: Sequence[T], f: Callable[[T], Any]):
    hparam = choose_single(hparam_options)
    return hparam, f(hparam)

# Handlers.
def choose_enumerate_handler[T](options: Sequence[T], k, lk):
    """Enumerative choice handler, using a vectorized loss function."""
    losses = jax.vmap(lk)(jnp.asarray(options))
    best_option = options[jnp.argmin(losses)]
    return k(best_option)

# Usage.
hyperparameter_tuning = handle(
    Handler(choose_single=choose_enumerate_handler),
)(hyperparameter_tuning)

lrs = [0.001, 0.002, 0.005, ...]
params, x, y = ...

def f(lr: float):
    g = handle(
        Handler(parameterized=True, choose=choose_grad_handler),
    )(linear_regression)
    return g(lr)(params, x, y)

loss_value, (best_hparam, new_params) = hyperparameter_tuning(lrs, f)
```

Figure 5: Hyperparameter tuning with CHOIX.

## D   Reinforcement learning

```python
# Effect operations for reinforcement learning.
@effect
def predict(state: State) -> Action:
    """Selects an action given an environment state."""

@effect
def observe(state: State, action: Action) -> tuple[Reward, State]:
    """Performs an action to get a reward and new environment state."""

Feedback = tuple[State, Action, Reward, State] # (old_state, action, reward, new_state)

@effect
def feedback(data: Feedback) -> None:
    """Updates agent with reward feedback."""

# Generic reinforcement learning algorithm.
def run_episode(state: State, is_goal: Callable[[State], bool]) -> State:
    """Runs reinforcement learning for one episode, updating environment with actions from agent."""
    def cond(state: State):
        return jnp.invert(is_goal(state))

    def body(state: State):
        action = predict(state)
        reward, new_state = observe(state, action)
        feedback((state, action, reward, new_state))
        return new_state

    return lax.while_loop(cond, body, state)

# Random number generation handler.
def random_uniform_handler(key: PRNGKey, args: tuple[int, int], k, lk):
    minval, maxval = args
    key, subkey = random.split(key)
    return k(key, random.uniform(subkey, minval, maxval))

# Agent handlers: Q-learning.
def predict_handler(params: Theta, state: State, k, lk):
    # Selects an action given agent parameters and state.
    epsilon = random_uniform(0, 1)
    action = jax.lax.select(
        epsilon < 0.1,
        jnp.floor(random_uniform(0, action_count)),
        find_max_action(params, state)
    )
    return k(params, action)

def feedback_handler(params: Theta, data: Feedback, k, lk):
    # Updates agent parameters given feedback.
    max_reward = find_max_reward(params, data)
    new_params = update(params, max_reward, data)
    return k(new_params)

# Environment handler: cliff-walking.
# Reference: R. Sutton and A. Barto, "Reinforcement Learning: An Introduction", 2020.
def cliff_walking_observe_handler(state: State, action: Action, k, lk):
    reward, new_state = cliff_walking_environment_update(state, action)
    return k(reward, new_state)

def cliff_walking_is_goal(state: State) -> bool:
    """Returns True if state is a goal state."""
    ...

# Usage.
run_qlearning = handle(
    Handler(parameterized=True, random_uniform=random_uniform_handler),
    Handler(parameterized=True, predict=predict_handler, feedback=feedback_handler),
    Handler(observe=cliff_walking_observe_handler),
)(run_episode)

loss_value, (new_rng_key, new_params, new_state) = run_qlearning(rng_key, params)(state, cliff_walking_is_goal)
```

Figure 6: A generic reinforcement learning algorithm in CHOIX (abbreviated for space).

The top-level `run_episode` function is generic and can be instantiated with different environments and agent algorithms through customized handlers. An outer "agent" handler predicts actions and applies feedback, while an inner "environment" handler updates environment state given actions.

Handlers here implement Q-learning (for `predict` and `feedback`) and cliff-walking (for `observe`). The Q-learning `predict_handler` itself uses a `random_uniform` effect for random number generation.

## E    Deep reinforcement learning

```python
# Reuses functions like `run_episode` from previous example.

@effect
def choose[T](x: T) -> T:
  """Choice effect operation."""

def atari_observe_handler(state: AtariState, action: Action, k, lk):
    reward, new_state = atari_environment_update(state, action)
    return k(reward, new_state)

def ppo_predict_handler(params: Theta, state: AtariState, k, lk):
    # Applies neural network to get per-action logits, then samples one.
    logits, value = forward.apply(params, state)
    ...
    pi = distrax.Categorical(logits)
    action = pi.sample()
    return k(params, action)

def ppo_feedback_handler(params: Theta, data: Feedback, k, lk):
    # Updates agent parameters given feedback.
    new_params = choose(params)
    # PPO loss has three parts: value loss, actor loss, entropy.
    # Intermediate computations are hidden.
    ...
    loss(value_loss)
    ...
    loss(actor_loss)
    ...
    loss(entropy)
    return k(new_params, None)

def choose_optax_handler(opt_state: Opt, params: Theta, k, lk):
    # Gradient-based update via the Optax library.
    grads = jax.grad(lk, argnums=1)(opt_state, params)
    updates, opt_state = optimizer.update(grads, opt_state)
    new_params = optax.apply_updates(params, updates)
    return k(opt_state, new_params)

run_ppo = handle(
    Handler(parameterized=True, random_uniform=random_uniform_handler),
    Handler(parameterized=True, predict=ppo_feedback_handler, feedback=ppo_feedback_handler),
    Handler(parameterized=True, choose=choose_optax_handler),
    Handler(observe=atari_observe_handler),
)(run_episode)

# Usage.
optimizer = optax.adam(learning_rate=...)
opt_state = optimizer.init(...)
loss_value, result = run_ppo(rng_key, params, opt_state)(state, cliff_walking_is_goal)
new_rng_key, new_params, new_opt_state, new_state = result
```

Figure 7: A deep reinforcement learning example in CHOIX.

This reuses the `run_episode` generic function (from Appendix D) with new handlers for `predict` and `feedback`. Here, handlers implement the PPO algorithm [8], abbreviated for simplicity.

In our CHOIX artifact, we use `run_ppo` to train agents on Atari game environments imported from the pgx library[8].

---

[8]https://github.com/sotetsuk/pgx