

# **Python Data structures**

# What we'll cover

- Lists
- Tuples
- Range
- Sets
- Dictionaries
- Strings
- Functions

A large yellow geometric shape, resembling a parallelogram or a trapezoid, is positioned in the top-left corner of the slide. It has a horizontal top edge, a horizontal bottom edge, and a slanted right edge that slopes downwards from left to right.

# Lists

# List basics

- Built-in collection of objects
- Mutable: its elements *can* be modified
- Arbitrarily typed elements (even differently typed)
  - Usually elements are of the same type
- Typically implemented as dynamic arrays (like `ArrayList` in Java)
- Constructors:
  - `[a, b, c]` — comma separated values, potentially empty
  - `list(iterable)` — e.g., `list(range(5))`
  - `[x for x in iterable]` — list comprehension, more later

# List comprehension

- Create lists with simple computations
- “... brackets containing an expression followed by a for clause, then zero or more for or if clauses”
  - <https://docs.python.org/3/tutorial/datastructures.html#list-comprehensions>

# Basic list operations

- Running example: `a = [3, 1, 4, 5]`
- `len(a)` — returns the number of elements in `a`: 4
- `a[idx]` — returns the `idx`-th element of `a`, using 0-based indexing
  - `idx` must be between 0 and `len(a) - 1`
  - `a[2] → 4`
- `a[-idx]` — returns the `idx`-th element of `a` from the right, using 1-based indexing
  - Intuitively means `len(a) - idx`, hence 1-based
  - `idx` must be between 1 and `len(a)`
  - `a[-1] → 5`

# List slicing

- Running example: `a = [3, 1, 4, 5]`
- Slices are **copies** of (not necessarily contiguous) subarrays
- `a[idx_0 : idx_f]` — returns a slice of `a` from `idx_0` (inclusive) to `idx_f` (non-inclusive). Both indices are optional.
  - Defaults: `idx_0 = 0, idx_f = len(a)`
  - `a[1:3] → [1, 4]`
  - `a[1:] → [1, 4, 5]`
  - `a[:3] → [3, 1, 4]` (equivalently `a[:-1]`)
  - `a[:] → [3, 1, 4, 5]` (copy)
- `a[idx_0 : idx_f : step]` — as above, but in steps of `step`
  - `a[0:3:2] → [3, 4]`
  - `a[-2:-5:-2] → [4, 3]`
  - `a[::-1] → [5, 4, 1, 3]`
- Support slice assignment: requires assigned array to be the same size as the slice
  - `a[0:3:2] = [0, 0] → [0, 1, 0, 5]`

# List implementation

<code>len(a)</code>	Number of elements	$O(1)$	
<code>a[idx_0:idx_f]</code>	Slice	$O(\text{idx}_f - \text{idx}_0)$	$O(1)$ for single element
<code>a + b</code>	Concatenate lists	$O(n+m)$	<code>+=</code> for concatenate update
<code>a * k</code>	Repeat list k times	$O(kn)$	<code>*=</code> for repeat update
<code>x in a</code>	True if some element in a has value x	$O(n)$	Linear search
<code>min(a),</code>	Minimum and maximum values	$O(n)$	
<code>max(a)</code>	Position of first element with value	$O(n)$	
<code>a.index(x)</code>	x		
<code>a.count(x)</code>	Number of elements with value x	$O(n)$	
<code>del a[idx_0:idx_f]</code>	Remove slice from a	$O(n)$	$O(n)$ even for single element
<code>a.sort()</code>	Sort list in place	$O(n \log n)$	
<code>a.pop(-idx)</code>	Delete the idx-th element	$O(\text{idx})$	$O(1)$ for last element (default)
<code>a.append(x)</code>	Add element with value x at the end	$O(1)$	



A large yellow geometric shape, resembling a parallelogram or a trapezoid, is positioned in the top-left corner of the slide. It has a horizontal top edge, a horizontal bottom edge, and a slanted right edge that slopes downwards from left to right.

# Tuples

# Tuple basics

- Built-in collection of objects
- Immutable: its elements are fixed after creation
- Arbitrarily typed elements (even differently typed)
  - Common to have multiple types in a tuple
  - E.g., a number and associated name ( ' a ' , 1 )
- Typically implemented as a static array (`Array` in Java)
- Constructors:
  - `()` — empty tuple
  - `(a,)` — singleton tuple
  - `(a,b,c)` or `a,b,c` — multiple elements
  - `tuple(iterable)` — created from iterable, order preserved

# Notes on tuples

- They are faster than lists
  - Use them if you won't need to modify it at runtime!
- There is no tuple comprehension
  - `(x for x in range(3))` just creates a *generator*, not a tuple
- Support all list operations for accessing, but not for altering

# Tuple implementation

Operation	Description	Runtime	Note
<code>len(a)</code>	Number of elements	$O(1)$	
<code>a[idx_0:idx_f]</code>	Slice	$O(\text{idx}_f - \text{idx}_0)$	$O(1)$ for single element
<code>a + b</code>	Concatenate tuples	$O(n+m)$	
<code>a * k</code>	Repeat tuple $k$ times	$O(kn)$	
<code>x in a</code>	True if some element in $a$ has value $x$	$O(n)$	
<code>min(a), max(a)</code>	Minimum and maximum values	$O(n)$	
<code>a.index(x)</code>	Position of first element with value $x$	$O(n)$	
<code>a.count(x)</code>	Number of elements with value $x$	$O(n)$	
<code>sorted(a)*</code>	Sort tuple <i>not</i> in place	$O(n \log n)$	

\*Returns sorted *list*. Compare to lists' `a.sort()` in place method

A large yellow geometric shape, resembling a parallelogram or a trapezoid, is positioned in the top-left corner of the slide. It has a horizontal top edge, a horizontal bottom edge, and a slanted right edge that slopes downwards from left to right.

Range

# Range basics

- Built-in collection of `int` objects
- Immutable
- Its size is  $O(1)$
- Constructors
  - `range(stop)` — from 0 to `stop-1`
  - `range(start, stop)` — from `start` to `stop-1`
  - `range(start, stop, step)` — from `start` to `stop-1` on steps of `step`
- Like tuples, support accessing operations
- Does *not* support the repetition (`*`) operator

# Sequences

- Types we've seen (`list`, `tuple`, `range`) are *sequence* types
- We will see another sequence type, `string`, later today
- Support for common indexing and slicing operations
- Sequence comparison is done in lexicographical order
  - Compare first element, if equal move to second, and so on
  - If an element is itself a sequence, compare recursively
  - Sequences must be of the same type (e.g., two `lists`)

A large yellow geometric shape, resembling a parallelogram or a trapezoid, is positioned in the top-left corner of the slide. It has a horizontal top edge, a horizontal bottom edge, and a slanted right edge that slopes downwards from left to right.

# Sets



# Set basics

- Mutable collection of objects with no repeated elements
- No ordering imposed
- Implemented as a hash set
- Supports only hashable types, but may contain multiple types
  - Mutable types are *not* hashable
- Constructors
  - `set()` — an empty set (cannot use `{ }`, reserved for dictionaries)
  - `set(iterable)` — non-repeated elements, order not preserved
  - `{a,b,c}` — set from elements, ignore repeated
  - `{x for x in iterable}` — set comprehension, just like lists

# Set implementation

Operation	Description	Runtime	Note
<code>len(a)</code>	Number of elements	$O(1)$	
<code>x in a</code>	True if some element in a has value x	$O(1)$	
<code>a.add(x)</code>	Insert element if not repeated	$O(1)$	Hashing
<code>a.remove(x)</code>	Remove element, error if not present	$O(1)$	
<code>a.discard(x)</code>	Remove element if present	$O(1)$	
<code>a &lt;= b</code>	True if a is improper subset of b	$O(n)$	< for proper
<code>a &gt;= b</code>	True if a is improper superset of b	$O(m)$	> for proper
<code>a.isdisjoint(b)</code>	True if no element of a is in b	$O(\min(n,m))$	
<code>a   b</code>	Union of sets	$O(n+m)$	= for union update
<code>a &amp; b</code>	Intersection of sets	$O(\min(n,m))$	&= for intersection update
<code>a - b</code>	Difference of sets	$O(n)$	-= for difference update
<code>a ^ b</code>	Symmetric difference of sets	$O(n)$	^= for sym. diff. update
<code>min(a), max(a)</code>	Minimum and maximum values	$O(n)$	Linear search
<code>sorted(a)*</code>	Sort set <i>not</i> in place	$O(n \log n)$	

\*Returns sorted *list*. Compare to lists' `a.sort()` in place method

A large yellow geometric shape, resembling a parallelogram or a trapezoid, is positioned in the top-left corner of the slide. It has a diagonal edge separating it from the white background.

# Dictionaries

# Dictionary basics

- Mutable collection of key-value pairs with no repeated elements
- As of Python 3.7, insertion order is preserved
- Implemented as a hash table
- Keys may only be hashable types, values are arbitrary types
- Constructors
  - `dict()`, `{}` — empty dictionary
  - `dict(iterable)` — sequence of k-v pairs
    - E.g., `dict([('a',1),('b',2)])`
  - `{k1: v1, k2: v2, k3: v3}` — key value pairs
  - `{k : v for k,v in iterable}` — dictionary comprehension
    - E.g., `{x : x**2 for x in range(5)}`

# Dictionary implementation

Operation	Description	Runtime	Note
<code>len(a)</code>	Number of elements	$O(1)$	
<code>key in a</code>	True if some element in a has key key	$O(1)$	Hashing
<code>a[key] = val</code>	Insert key-val pair	$O(1)$	Hashing, overwrite if present
<code>del a[key]</code>	Remove element, error if not present	$O(1)$	Hashing
<code>a[key]</code>	Return val, error if key not present	$O(1)$	Hashing
<code>a.update(b)</code>	Update with key-val from b	$O(m)$	Overwrite if present
<code>a.keys()</code>	Return <i>view</i> of keys	$O(1)$	Supports set operations
<code>a.values()</code>	Return <i>view</i> of values	$O(1)$	
<code>a.items()</code>	Return view of (key, val) pairs	$O(1)$	

- collections.Counter subclasses dictionaries. Useful to count instances
  - <https://docs.python.org/3/library/collections.html#collections.Counter>

# Looping through dictionaries

- If you only need the keys, use dictionary as iterable
- If you need both keys and values, use `items ( )`

<https://docs.python.org/3/tutorial/datastructures.html#looping-techniques>

A solid yellow geometric shape, resembling a parallelogram or a trapezoid, is positioned in the top-left corner of the slide. It has a vertical left edge, a horizontal top edge, and a diagonal right edge that slopes downwards from left to right.

# Strings

# String basics

- Python's built-in text sequence type
- Stored as Unicode
- Immutable sequence: supports sequence accessing operations
- Comparisons are lexicographical



# String constructors

- `'We can use "double" quotes'` — single quotations
- `"We can use 'single' quotes"` — double quotations
- `'''Triple'''` or `"""Triple"""` — triple quotes
- `str(obj)` — uses obj's `__str()` method
  - Does *not* parse iterables with characters or strings
  - E.g., `str(['a', 'b', 'c'])` → `"['a', 'b', 'c']"`
  - More on this later in the course
- `str()` or `''` or `"""` — empty string

# String concatenation

- String literals are concatenated with whitespaces
  - `'Hello' 'there' → "Hello there"`
- String variables are concatenated with `+`
- Repeat string with `*`
- Concatenate strings from iterable with `sep.join(iterable)`
- Immutable type! No support for `+=` or `*=`

# String indexing

- Strings can be indexed like sequences
- Slicing works as with sequences too
- Indexing and slicing return substrings
  - E.g., `'Python'[1:3]`  $\rightarrow$  `'yt'`
- There is no special character type
  - Characters are just strings of length 1

# String matching

- `b in a` — whether `b` is a substring of `a`
- `a.find(b, idx_0, idx_f)` — position of first occurrence of `b` in `a`
  - If `b` is not in `a`, returns `-1`
  - `a.index()` has the same notation but raises error if `b` not in `a`
  - `a.rfind()` or `a.rindex()` return the position of last occurrence
- `a.count(b, idx_0, idx_f)` — occurrences of `b` in `a`
- `a.startswith(prefix, idx_0, idx_f)` — whether `a` begins with `prefix`
  - `prefix` may be a tuple of strings
  - `a.endswith()` for suffixes
- `idx_0` and `idx_f` are optional, and interpreted in slice notation
  - Why is this better than simply slicing?

# String formatting

- `a.format(x, y, name=z, ...)`
- `a` must contain *replacement fields*
  - Expressions surrounded by `{ }`
  - To include `{, }` in `a`, escape them as double `{{, }}`
- `{idx!conversion:format_spec}`  
`{name!conversion:format_spec}`
  - `idx` — number of the argument (optional, but must be consistent)
  - `name` — may use this for named arguments (possibly mixed with `idx`)
  - If neither is provided, go in order
  - `conversion` — optional, `!s` (`str()`), `!a` (`ascii()`) or `!r` (`repr()`)
  - `format_spec` — optional, details to follow

# Other string operations

Operation	Description	Note	Related
<code>len(a)</code>	Number of characters in <code>a</code>		
<code>a.capitalize()</code>	First character capitalized, rest lowercased		<code>lower()</code> , <code>swapcase()</code> , <code>title()</code>
<code>a.expandtabs(t)</code>	Replace tabs with spaces to fill tabs of size <code>t</code>		
<code>a.isalnum()</code>	True if all characters are alphanumeric		<code>isalpha()</code> , <code>isdecimal()</code> , <code>isdigit()</code> , <code>isnumeric()</code>
<code>a.islower()</code>	True if all characters are lowercase		<code>istitle()</code> , <code>isupper()</code>
<code>a.isspace()</code>	True if all characters are whitespace		
<code>a.partition(sep)</code> <code>a.replace(old, new)</code>	Split string into prefix, <code>sep</code> , and suffix Replace all occurrences of <code>old</code> by <code>new</code>	If <code>sep</code> is not found, return <code>a</code> and two empty strings	<code>rpartition()</code>
<code>a.split(sep, maxsplit)</code>	Return list of words with <code>sep</code> as delimiter. Default: whitespace	At most <code>maxsplit</code> splits (optional), leftmost	<code>rsplit()</code>
<code>a.splitlines()</code>	Return list of lines, removing line breaks		
<code>a.strip()</code>	Remove leading and trailing spaces	Optionally: pass string to remove specific characters	<code>rstrip()</code> , <code>lstrip()</code>

Note: all functions for modifying the string return a modified copy

A large yellow geometric shape, resembling a parallelogram or a trapezoid, is positioned in the top-left corner of the slide. It has a horizontal top edge, a horizontal bottom edge, and a slanted right edge that slopes downwards from left to right.

# Functions

# Function definitions

- Prototypical example

```
def function(x, y, z=0):  
    print(x, y, z)
```

`function(1,2) → 1, 2, 0`

- Only optional parameters may come after the first optional parameter
- Default parameters are evaluated at the time of definition
  - This is crucial if you want to use a mutable type!