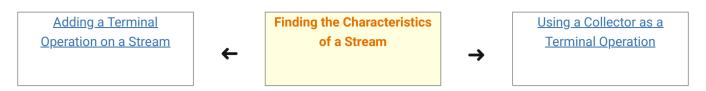
Home > Tutorials > The Stream API > Finding the Characteristics of a Stream



Finding the Characteristics of a Stream

Characteristics of a Stream

The Stream API relies on a special object, an instance of the <u>Spliterator</u> interface. The name of this interface comes from the fact that the role of a spliterator in the Stream API looks like the role of an iterator has in the Collection API. Moreover, because the Stream API supports parallel processing, a spliterator object also controls how a stream splits its elements among the different CPUs that handle parallelization. The name is the contraction of *split* and *iterator*.

Covering this spliterator object in details is beyond the scope of this tutorial. What you need to know is that this spliterator object carries the *characteristics* of a stream. These characteristics are not something you will often use, but knowing what they are will help you to write better and more efficient pipelines in certain cases.

The characteristics of a stream are the following.

Characteristic	Comment
<u>ORDERED</u>	The order in which the elements of the stream are processed matters.
DISTINCT	There are no doubles in the elements processed by that stream.
NONNULL	There are no null elements in that stream.
SORTED	The elements of that stream are sorted.
SIZED	The number of elements this stream processes is known.
SUBSIZED	Splitting this stream produces two <u>SIZED</u> streams.

There are two characteristics, <u>IMMUTABLE</u> and <u>CONCURRENT</u>, which are not covered in this tutorial.

Every stream has all these characteristics set or unset when it is created.

Remember that a stream can be created in two ways.

- 1. You can create a stream from a source of data, and we covered several different patterns.
- 2. Every time you call an intermediate operation on an existing stream, you create a new stream.

The characteristics of a given stream depend on the source it has been created on, or the characteristics of the stream with which it was created, and the operation that created it. If your stream is created with a source, then its characteristics depend on that source, and if you created it with another stream, then they will depend on this other stream and the type of operation you are using.

Let us present each characteristic in more details.

Ordered Streams

<u>ORDERED</u> streams are created with ordered sources of data. The fist example that may come to mind is any instance of the <u>List</u> interface. There are others: <u>Files.lines(path)</u> and <u>Pattern.splitAsStream(string)</u> also produce <u>ORDERED</u> streams.

Sorted Streams

A <u>SORTED</u> stream is a stream that has been sorted. This stream can be created from a sorted source, such as a <u>TreeSet</u> instance, or by a call to the <u>sorted()</u> method. Knowing that a stream has already been sorted may be used by the stream implementation to avoid sorting again an already sorted stream. This optimization may not be used all the time because a <u>SORTED</u> stream may be sorted again with a different comparator than the one used the first time.

There are some intermediate operations that clear the <u>SORTED</u> characteristic. In the following code, you can see that both <u>strings</u> and <u>filteredStream</u> are <u>SORTED</u> streams, whereas <u>lengths</u> is not.

```
Collection<String> stringCollection = List.of("one", "two", "two", "three", "four", "five");

Stream<String> strings = stringCollection.stream().sorted();

Stream<String> filteredStrings = strings.filtered(s -> s.length() < 5);

Stream<Integer> lengths = filteredStrings.map(String::length);
```

Mapping or flatmapping a <u>SORTED</u> stream removes this characteristic from the resulting stream.

Distinct Streams

A <u>DISTINCT</u> stream is a stream with no duplicates among the elements it is processing. Such a characteristic is acquired when building a stream from a <u>HashSet</u> for instance, or from a call to the

distinct() intermediate method call.

The <u>DISTINCT</u> characteristic is kept when filtering a stream but is lost when mapping or flatmapping a stream.

Let us examine the following example.

```
Collection<String> stringCollection = List.of("one", "two", "two", "three", "four", "five");

Stream<String> strings = stringCollection.stream().distinct();

Stream<String> filteredStrings = strings.filtered(s -> s.length() < 5);

Stream<Integer> lengths = filteredStrings.map(String::length);
```

- <u>stringCollection.stream()</u> is not <u>DISTINCT</u> as it is build from an instance of <u>List</u>.
- strings is <u>DISTINCT</u> as this stream is created by a call to the <u>distinct()</u> intermediate method.
- filteredStrings is still <u>DISTINCT</u>: removing elements from a stream cannot create duplicates.
- length has been mapped, so the <u>DISTINCT</u> characteristic is lost.

Non-Null Streams

A <u>NONNULL</u> stream is a stream that does not contain null values. There are structures from the Collection Framework that do not accept null values, including <u>ArrayDeque</u> and the concurrent structures like <u>ArrayBlockingQueue</u>, <u>ConcurrentSkipListSet</u>, and the concurrent set returned by a call to <u>ConcurrentHashMap.newKeySet()</u>. Streams created with <u>Files.lines(path)</u> and <u>Pattern.splitAsStream(line)</u> are also <u>NONNULL</u> streams.

As for the previous characteristics, some intermediate operations can produce a stream with different characteristics.

- Filtering or sorting a *NONNULL* stream returns a *NONNULL* stream.
- Calling <u>distinct()</u> on a <u>NONNULL</u> stream also returns a <u>NONNULL</u> stream.
- Mapping or flatmapping a <u>NONNULL</u> stream returns a stream without this characteristic.

Sized and Subsized Streams

Sized Streams

This last characteristic is very important when you want to use parallel streams. Parallel streams are covered in more detail later in this tutorial.

A <u>SIZED</u> stream is a stream that knows how many elements it will process. A stream created from any instance of <u>Collection</u> is such a stream because the <u>Collection</u> interface has a <u>size()</u> method, so getting

this number is easy.

On the other hand, there are cases where you know that your stream will process a finite number of elements, but you cannot know this number unless you process the stream itself.

This is the case for streams created with the <u>Files.lines(path)</u> pattern. You can get the size of the text file in bytes, but this information does not tell you how many lines this text file has. You need to analyze the file to get this information.

This is also the case for the <u>Pattern.splitAsStream(line)</u> pattern. Knowing the number of characters there are in the string you are analyzing does not give any hint about how many elements this pattern will produce.

Subsized Streams

The <u>SUBSIZED</u> characteristic has to do with the way a stream is split when computed as a parallel stream. In a nutshell, the parallelization mechanism splits a stream in two parts and distribute the computation among the different available cores on which the CPU is executing. This splitting is implemented by the instance of the <u>Spliterator</u> the stream uses. This implementation depends on the source of data you are using.

Suppose that you need to open a stream on an <u>ArrayList</u>. All the data of this list is held in the internal array of your <u>ArrayList</u> instance. Maybe you remember that the internal array on an <u>ArrayList</u> object is a compact array because when you remove an element from this array, all the following elements are moved one cell to the left so that no hole is left.

This makes the splitting an <u>ArrayList</u> straightforward. To split an instance of <u>ArrayList</u>, you can just split this internal array in two parts, with the same amount of elements in both parts. This makes a stream created on an instance of <u>ArrayList</u> <u>SUBSIZED</u>: you can tell in advance how many elements will be held in each part after splitting.

Suppose now that you need to open a stream on an instance of HashSet stores its elements in an array, but this array is used differently than the one used by ArrayList. In fact, more than one element can be stored in a given cell of this array. There is no problem in splitting this array, but you cannot tell in advance how many elements will be held in each part without counting them. Even if you split this array by the middle, you can never be sure that you will have the same number of elements in both halves. This is the reason why a stream created on an instance of HashSet is SIZED but not SUBSIZED.

Transforming a stream may change the SIZED and SUBSIZED characteristics of the returned stream.

- Mapping and sorting a stream preserves the <u>SIZED</u> and <u>SUBSIZED</u> characteristics.
- Flatmapping, filtering, and calling distinct() erases these characteristics.

It is always better to have <u>SIZED</u> and <u>SUBSIZED</u> stream for parallel computations.

Last update: September 14, 2021

Adding a Terminal
Operation on a Stream



Finding the Characteristics of a Stream



<u>Using a Collector as a</u> <u>Terminal Operation</u>

<u>Home</u> > <u>Tutorials</u> > <u>The Stream API</u> > Finding the Characteristics of a Stream