

## Ruby master - Feature #17001

### [Feature] Dir.scan to yield dirent for efficient and composable recursive directory scanning

06/30/2020 07:57 PM - byroot (Jean Boussier)

<b>Status:</b>	Open
<b>Priority:</b>	Normal
<b>Assignee:</b>	
<b>Target version:</b>	
<b>Description</b>	
<b>Use case</b>	
<p>When you need to recursively scan a directory, you either have to use <code>Dir[] / Dir.glob</code>, which is fine for small directories or simple patterns, but can easily take several seconds to complete for large repositories or complex patterns and returns a very large array which tend to trash GC.</p> <p>Or you can use <code>Dir.each_entry / Dir.foreach</code> recursively, but then you need to stat each entry to know whether it's a directory, or even symlink if you want to follow them.</p> <p>This means one syscall per directory, and one per file and directories. This is particularly impactful on OSX where <code>stat()</code> is several times slower than on Linux because of various sandboxing features.</p> <p>There's a <a href="#">typical example of this use case in Bootsnap</a>.</p>	
<b>Proposal</b>	
<p><a href="#">Python introduced os.scandir a few years ago</a> for exactly this purpose. It is functionally similar to <code>Dir.foreach / Dir.each_child</code>, except it yields <code>DirEntry</code> instances which are a wrapper around the libc <code>dirent</code> struct.</p> <p>I reduced the Bootsnap code into a <a href="#">simplified benchmark</a>, and using <code>os.scandir()</code> Python scan our main repo in a bit over 1s, which 3 to 4 times faster than Ruby can with <code>Dir.foreach</code> (3-4s). For comparison sake <code>Dir['**/*.rb']</code> also complete in about 1s.</p> <p>So I believe that exposing a similar <code>Dir.scan</code> method, returning <code>Dir::Entry</code> instances, with methods inspired from <code>File::Stat</code> such as <code>directory?</code> would allow for more performant file system scanning when the query is not easily expressed with a glob pattern.</p>	