(19) **United States**

(12) **Patent Application Publication** (10) Pub. No.: **US 2023/0232195 A1**

Biswas et al. (43) **Pub. Date:** **Jul. 20, 2023**

(54) **COLLECTIVE SCALING OF APPLICATIONS**

(71) Applicant: **VMware, Inc.**, Palo Alto, CA (US)

(72) Inventors: **Sudipta Biswas**, Bangalore (IN);
**Monotosh Das**, Bangalore (IN);
**Hemant Kumar Shaw**, Bangalore (IN);
**Shubham Chauhan**, Faridabad (IN)

**Publication Classification**

(57) **ABSTRACT**

Some embodiments provide a method for scaling a service chain that includes multiple services, each of which is provided by one or more instances of the service. The method identifies that a first service in the service chain has received a number of requests. For each service in the service chain, the method (i) identifies a scaling factor that estimates a portion of requests received at the first service that will be subsequently received at the service and (ii) deploys a number of additional instances of the service based on the identified scaling factor for the service and the number of requests received at the first service.