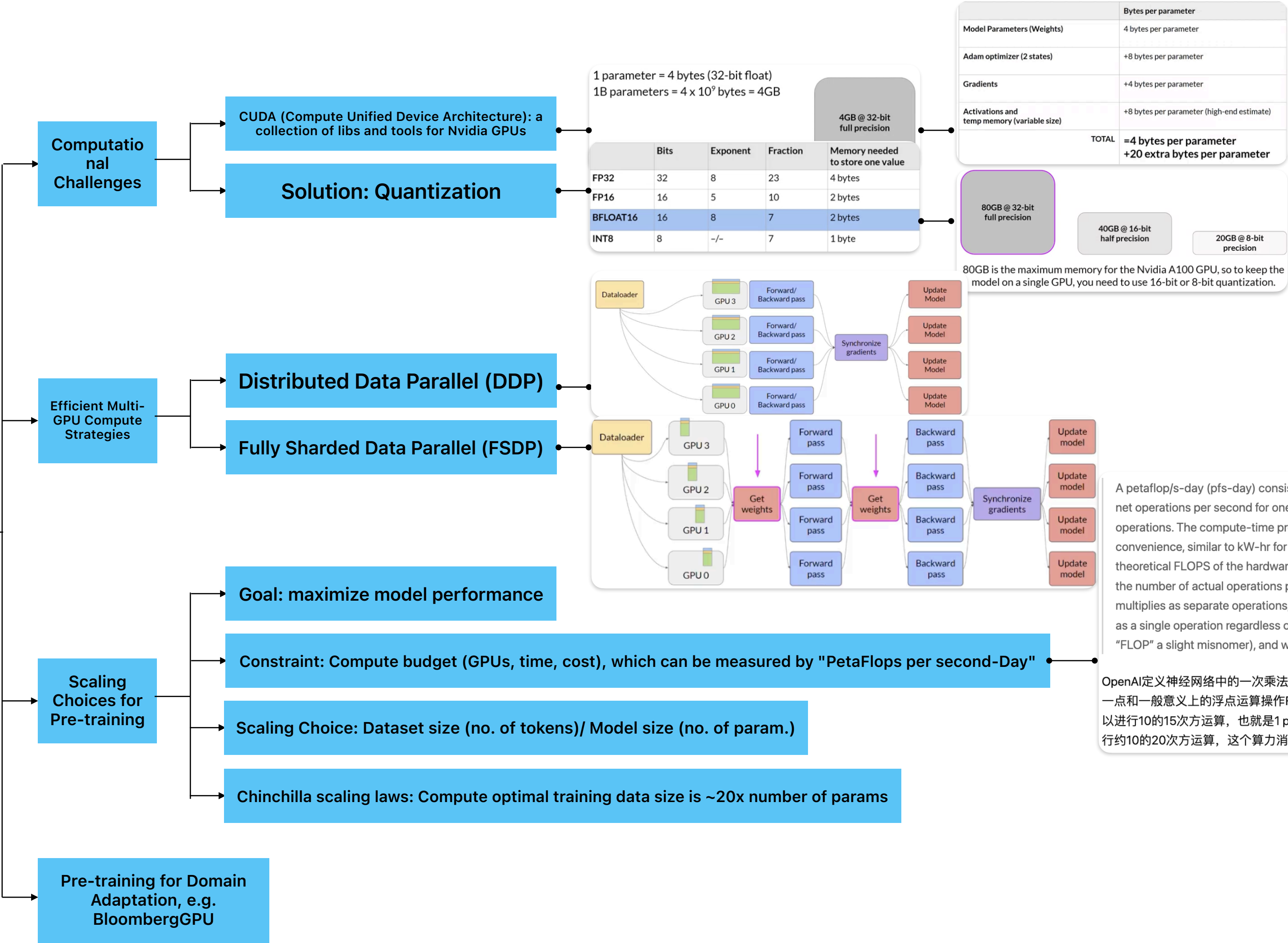


The Usage of LLMs



A petaflop/s-day (pfs-day) consists of performing  $10^{15}$  neural net operations per second for one day, or a total of about  $10^{20}$  operations. The compute-time product serves as a mental convenience, similar to kW-hr for energy. We don't measure peak theoretical FLOPS of the hardware but instead try to estimate the number of actual operations performed. We count adds and multiplies as separate operations, we count any add or multiply as a single operation regardless of numerical precision (making "FLOP" a slight misnomer), and we ignore ensemble models.

OpenAI定义神经网络中的一次乘法或者一次加法为一个操作，这一点和一般意义上的浮点运算操作FLOP略有不同。如果每秒钟可以进行 $10^{15}$ 次方运算，也就是1 peta flops，那么一天就可以进行约 $10^{20}$ 次方运算，这个算力消耗被称为1个petaflop/s-day。