

## PrimoCache

# Intel Optane SSD 900P 和 PrimoCache 详细性能评测及 在服务器/工作站的应用

文档编号	: RS-SW-PCC-00-18-01	版本	: 1.3
创建日期	: 2018-03-09	更新日期	: 2019-04-20
状态	: 发布	密级	: 公开

# 概述

Intel 基于 3D XPoint 技术的 Optane（中文名傲腾）系列产品凭借其超高性能、超低延迟和超长寿命，自推出以来便备受瞩目。迄今为止已经发布的产品包括面向数据中心的高端 SSD P4800X 系列、面向桌面缓存加速的傲腾内存系列，面向高端桌面、工作站和存储加速的 SSD 900P 系列和 800P 系列。

根据 Intel 官方发布的规格说明，SSD 900P 系列提供了半高半长 PCIe 卡和 U.2 2.5 英寸盘两种样式，目前支持的容量为 280GB 和 480GB，后续还会增加更高容量的型号。系统接口为 PCIe3.0 x4，标称性能最高顺序读写分别为 2500MB/s、2000MB/s，最高随机读写分别为 55 万 IOPS、50 万 IOPS，读写延迟时间低至 10  $\mu$ s。写入寿命为 18.69TB 每 GB 容量，按 5 年质保计算，相当于每天可 10 次全盘写入。

尽管 Intel 推出 SSD 900P 系列官方定位目标是消费级市场，但是 SSD 900P 在读写性能、延迟和 IO 吞吐量方面的表现非常杰出，可适合于各种工作站和中小型服务器。特别是当搭配 PrimoCache 软件使用时，可以将 SSD 900P 用作其它大容量低速硬盘的缓存，极大提升这些硬盘的读写速度和 IO 处理能力，甚至可以接近或达到 SSD 900P 的性能。对于没有足够预算采购昂贵的大容量高性能 SSD 的用户，PrimoCache + SSD 900P 缓存 + 大容量机械硬盘/低速 SSD 是一个极具性价比的方案。另一方面，这个缓存方案对工作站或服务器现有的硬件配置和软件环境无需做任何改动，也无需进行任何的数据迁移，只需要安装 SSD 900P 和 PrimoCache 软件即可带来超高的性能提升，相当简单方便，非常适合需要提升性能但不能改动已有硬件或软件的用户。

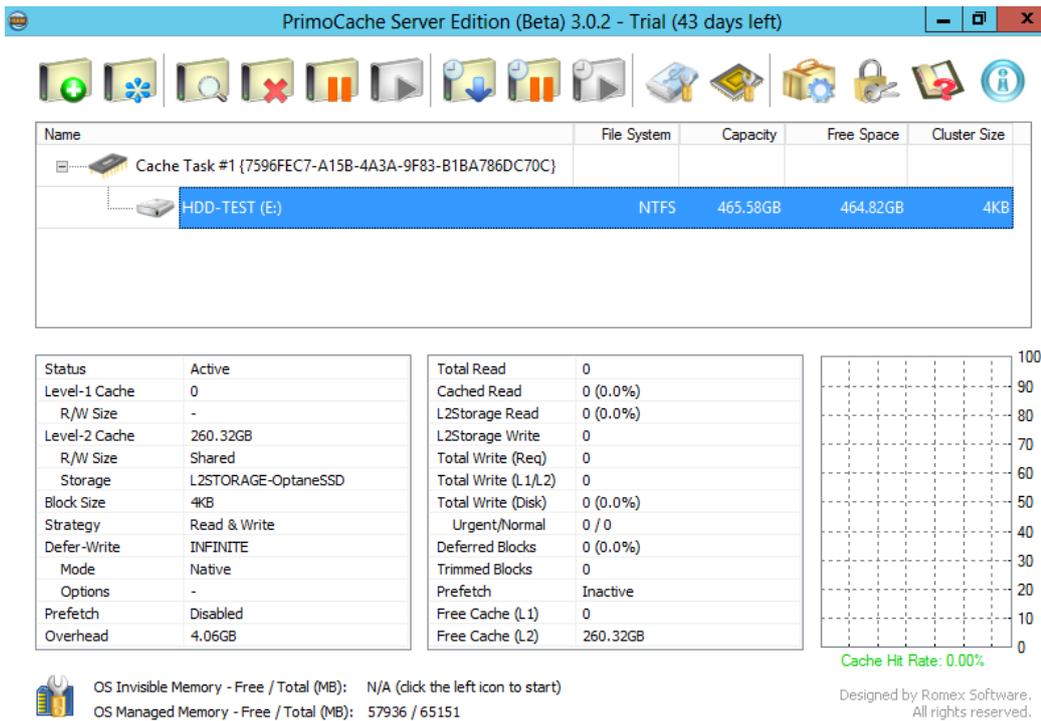
本文将详细测评 SSD 900P 自身的读写性能以及作为 PrimoCache 缓存的读写性能。本文不介绍如何使用 PrimoCache 软件以及如何将 SSD 900P 配置成缓存，这方面的说明请参见 PrimoCache [帮助文档](#) 或 [快速使用指南](#)。PrimoCache 软件可从 [官网下载页面](#) 下载试用。

# 测试方案

本文测试使用的 Optane 产品为 Intel Optane SSD 900P 280GB PCI-E 扩展卡。测试平台和系统软件环境如下表所示。

主板	Intel S2600CW
CPU	Intel Xeon CPU E5-2698 v3 @ 2.30GHz, 1 颗
内存	Samsung 16GB DDR4-2133 RDIMM x4, 共 64GB
硬盘	Seagate Constellation ES ST1000NM0011 1TB 3.5" Drive (7200RPM, SATA3, 64MB)
	Intel Optane SSD 900P 280GB PCI-E HDDL
操作系统	Windows Server 2012 Standard v6.2.9200
NVMe 驱动	Intel NVMe Driver 3.2.0.1002
测试工具	Microsoft Diskspd v2.0.17
PrimoCache	Server Edition v3.0.2 Beta

测试时使用的 PrimoCache 缓存配置如下图所示。目标硬盘仅设置了 SSD 缓存，没有一级内存缓存。



测试项目包括顺序标准读写、随机标准读写、顺序混合读写和随机混合读写，测试数据块大小覆盖 4KB 至 1MB 全系列。测试负载考虑了低中高各种情形，包括单线程单队深、单线程多队深、多线程单队深和多线程多队深，最高至 32 线程 64 队深。测试场景则涵盖满盘和空盘测试、全盘和区间（前后）测试以及寿命影响测试。整个测试方案基本上全方面的反映了测试对象在实际应用环境中各种可能情景的读写速度和 IO 性能。

## 测试报告格式说明

本文所有测试结果报告均以图表形式显示每组测试的测量结果。每份报告中每一行包含四个图表，依次显示了一组测试的四个测量值：数据传输率（MB/s）、IOPS、平均延时（毫秒）和平均 CPU 占用率（%）。

在标准读写测试和混合读写测试的测试结果报告中，每个图表通常会显示三条折线，分别标以“OptaneSSD 900P”、“ST1000NM0011”和“PrimoCache L2”。“OptaneSSD 900P”和“ST1000NM0011”分别指没有安装 PrimoCache 软件下对 Intel Optane SSD 900P 280GB 和 Seagate ST1000NM0011 硬盘测试的结果。“PrimoCache L2”指使用 SSD 900P 作为 ST1000NM0011 硬盘的缓存后，对 ST1000NM0011 硬盘进行测试的结果。

测试报告中的 QxTy 表示测试负载为 y 个线程同时操作，每个线程队列深度为 x，即同时 IO 操作数为 x\*y 个。比如 Q1T1 表示 1 线程 1 队深，单 IO 操作，Q64T32 表示 32 线程 64 队深，同时 IO 操作数目达 2048 个。

因测试结果较多，本文正文中仅引用了一小部分典型图表进行分析说明。如需查看全部详细数据和图表，请下载文末附件。

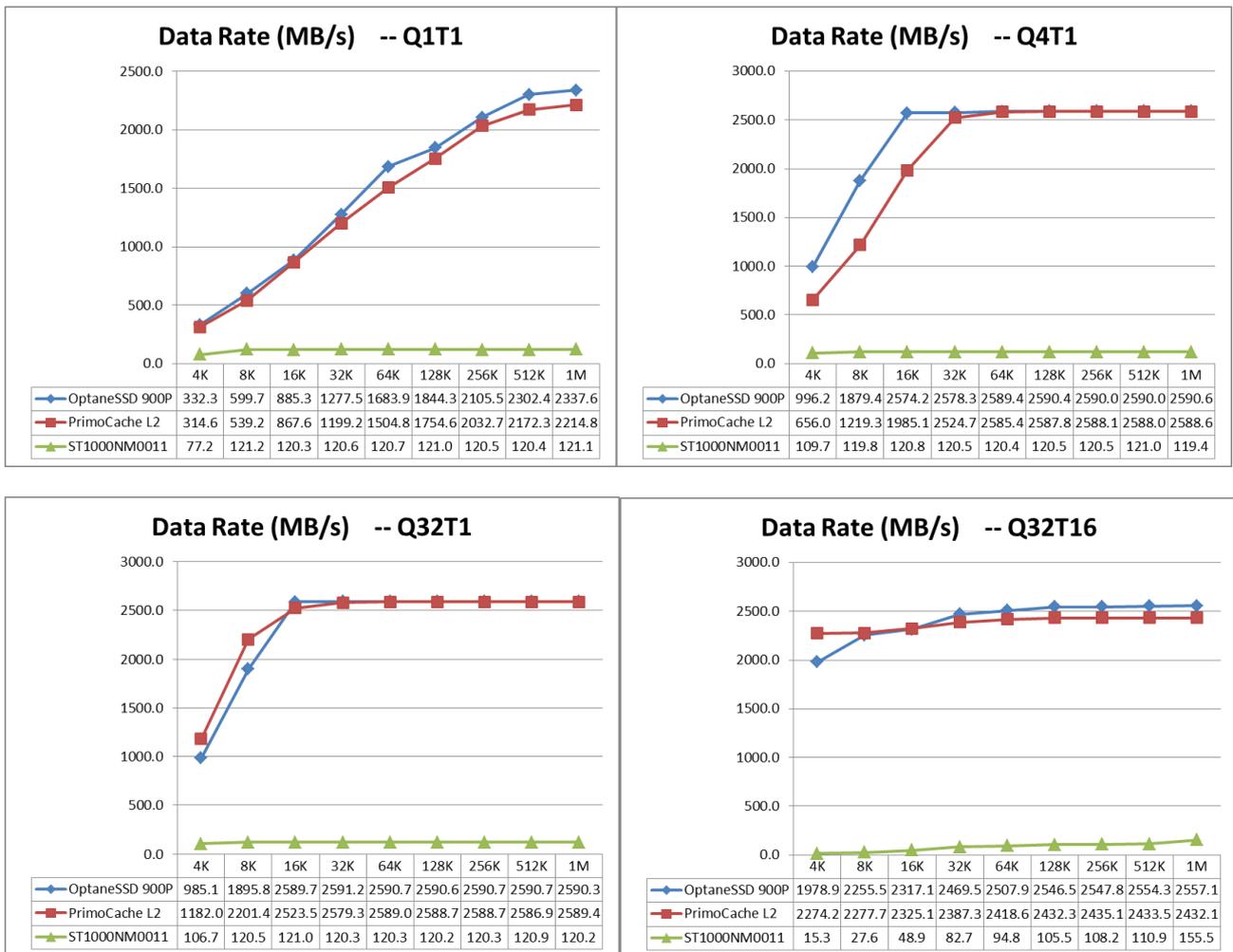
# 标准读写测试

标准读写测试主要测试顺序读、顺序写、随机读和随机写的性能。以下图表来自全盘测试的结果。全盘测试是指对整个存储设备的全部空间范围进行读写测试，在测试开始前测试文件的数据已经写满全盘，因此这里的全盘测试同时也是一个满盘测试。

## 顺序标准读

顺序标准读性能通常主要关注 64KB 至 1MB 大小数据块的处理能力。从以下图表可以看到，SSD 900P 在单线程单队深（极低负载）下的顺序读速度大约在 1700MB/s（64KB）至 2300MB/s（1MB），在多队深或多线程下其 64KB 至 1MB 的顺序读速度基本都稳定在官方标称 2500MB/s 以上，接近 2600MB/s！在服务器/工作站实际应用中，存储设备大多运行在多线程多队深状态下，因此 SSD 900P 可以充分发挥其最大性能。

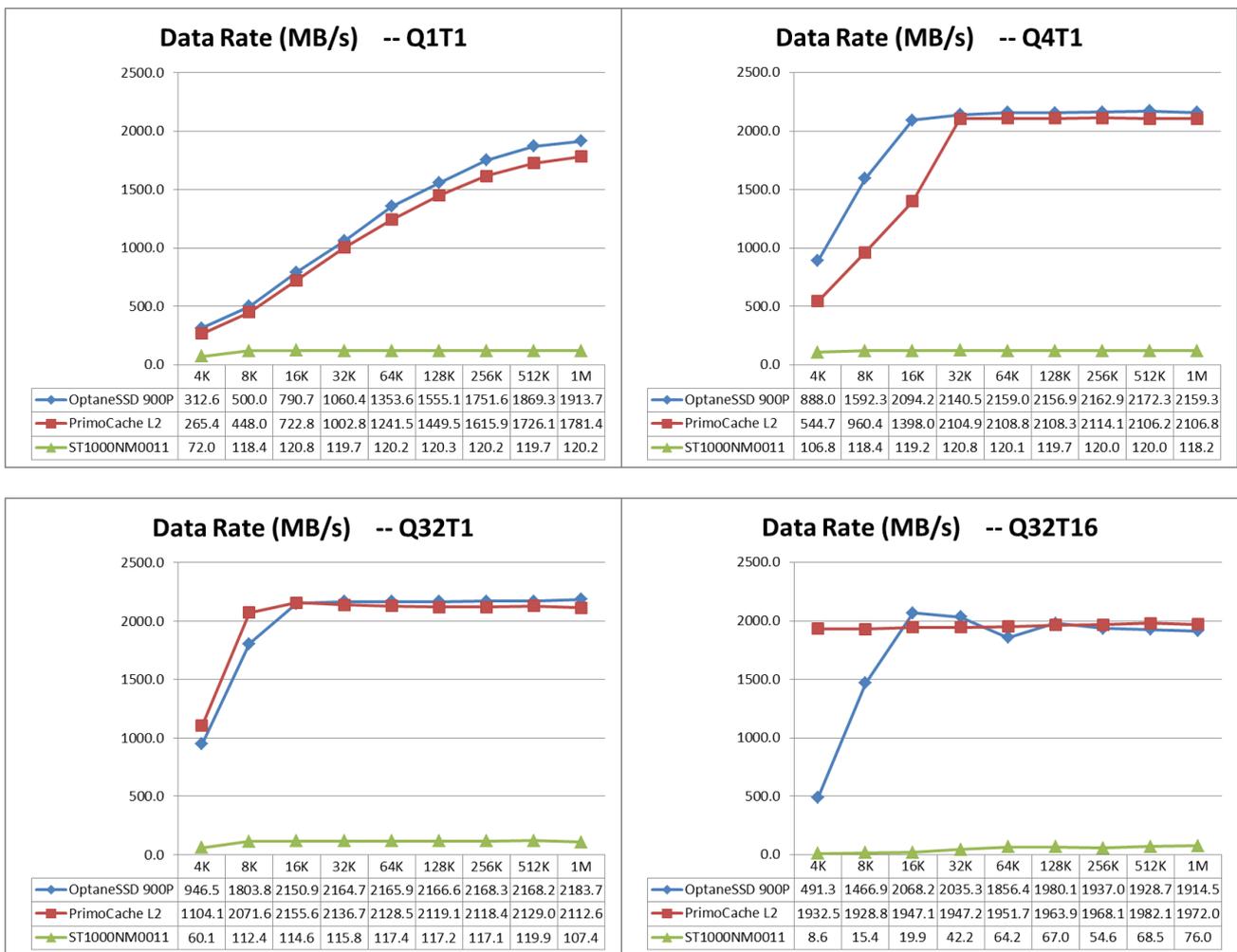
同时也可以看到，用 SSD 900P 作为 ST1000NM0011 硬盘的缓存后，ST1000NM0011 硬盘的顺序读速度接近或达到 SSD 900P 的性能，甚至在某些情况下，PrimoCache 缓存还可以充分挖掘 SSD 900P 的性能，使被缓存硬盘的性能超出 SSD 900P 自身。

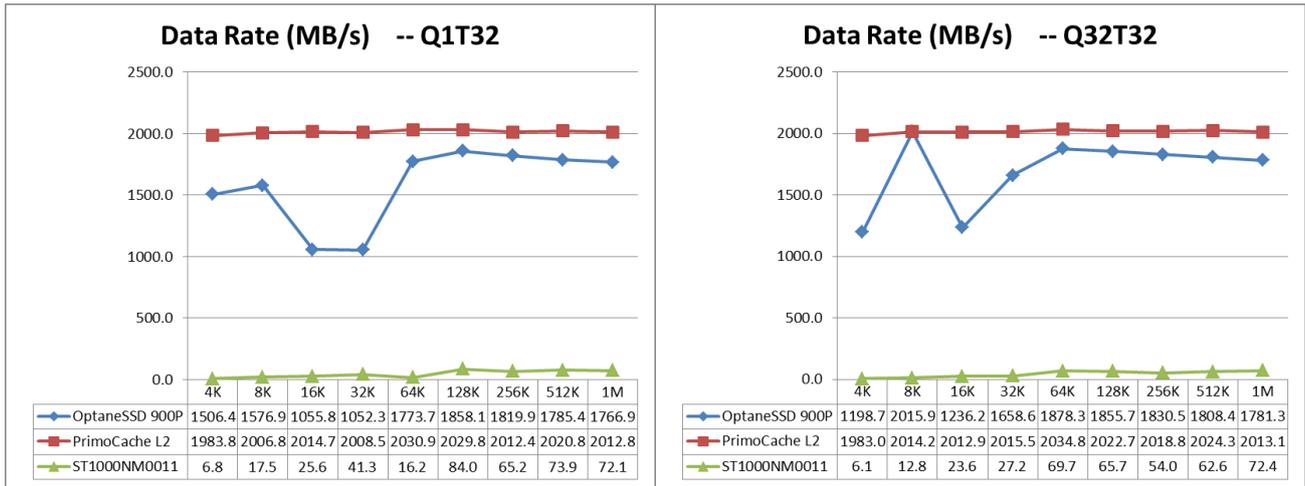


## 顺序标准写

顺序标准写性能，和顺序标准读类似，也主要关注 64KB 至 1MB 大小数据块的处理能力。从图表可以看到 SSD 900P 在单线程单队深下的顺序写速度大约在 1400MB/s（64KB）至 1900MB/s（1MB），在多队深或多线程（少于 32 线程）下其 64KB 至 1MB 的顺序写速度基本都稳定在官方标称 2000MB/s 上下，最高甚至接近 2200MB/s。

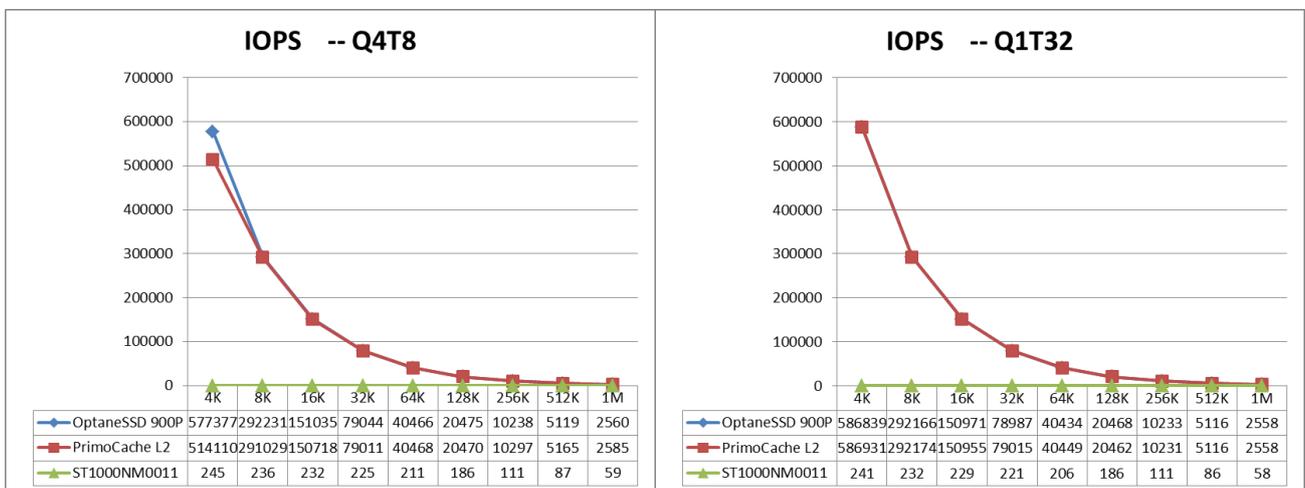
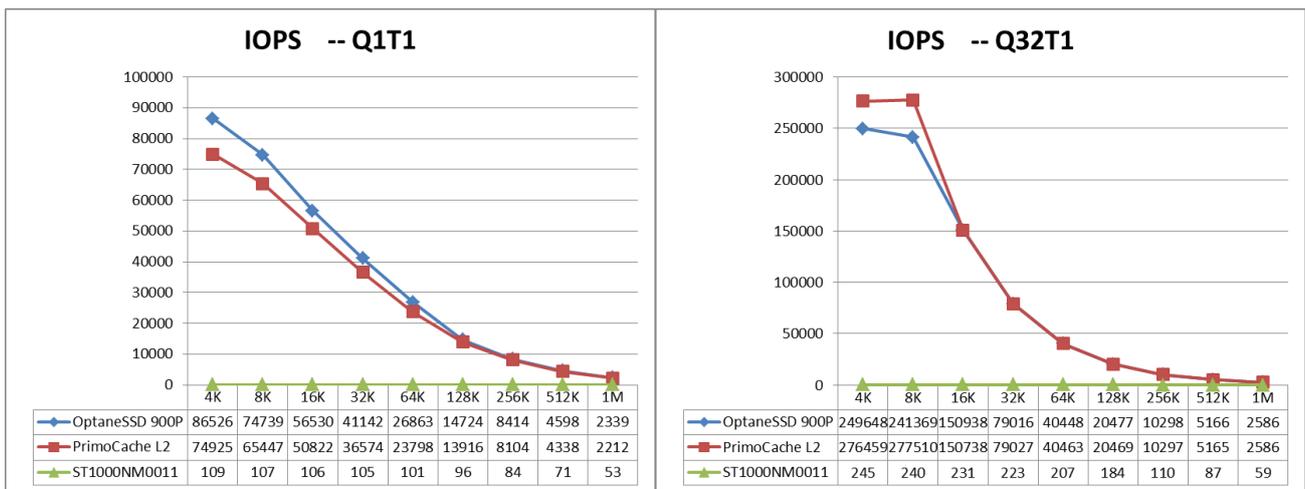
同样地，PrimoCache 缓存的顺序写性能也是基本接近或达到 SSD 900P 的性能。此外测试结果也显示了一个有意思的现象，随着同时操作的线程数增加，SSD 900P 的最高顺序写速度从大约 2200MB/s（单线程）降低到大约 1800MB/s（32 线程），然而 PrimoCache 缓存即使在 32 线程下仍然可以保持在 2000MB/s，最大程度的利用了 SSD 900P 的性能，使其性能反超 SSD 900P。猜测出现这种现象的一个可能原因是 SSD 900P 主控芯片受限于自身有限的硬件资源无法充分发挥存储介质的最大性能，而 PrimoCache 缓存则可以调配整个计算机系统的硬件资源达到最大性能。





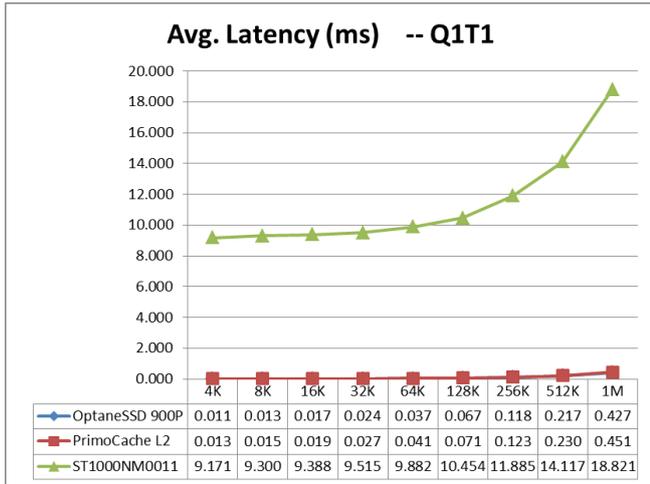
## 随机标准读

对于随机标准读，一般主要关注 4KB 至 64KB 大小数据块的处理性能，这里重点关注 4KB 数据块的 IOPS 数值。从图表中可以看到，SSD 900P 在单线程单队深下的 4KB 随机读 IOPS 数值接近 9 万，随着线程数或队深数的增加，其 IOPS 数值也是显著上升，在 8 线程 4 队深下，IOPS 已经接近 58 万。测试结果表明其 4KB 随机读 IOPS 最高数值稳定在近 59 万，比官方标称 55 万还要高一些。



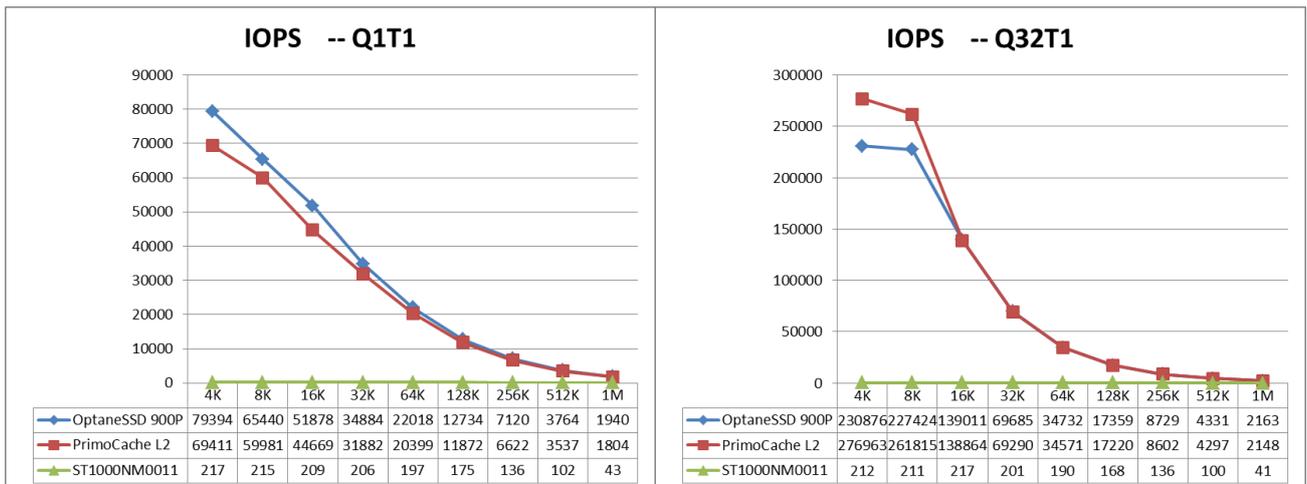
PrimoCache 缓存对于 4KB 数据块的处理性能在多数情形下会弱于 SSD 900P 本身。这是因为缓存处理 IO 时相对于存储设备直接处理 IO 会多一些操作时间，而 SSD 900P 对于 4KB 数据块的处理时间非常小，因此缓存额外处理的时间开销相对于整个 4KB 数据块处理时间的占比就会较大，反映在 IOPS 上的性能差异就比较明显，尤其在低线程数低队深的情况下。但是在低线程数高队深的情况下（比如 Q32T1），PrimoCache 缓存由于可以充分发挥 SSD 900P 的性能，其 4KB 数据块的 IOPS 值反而要高于 SSD 900P。

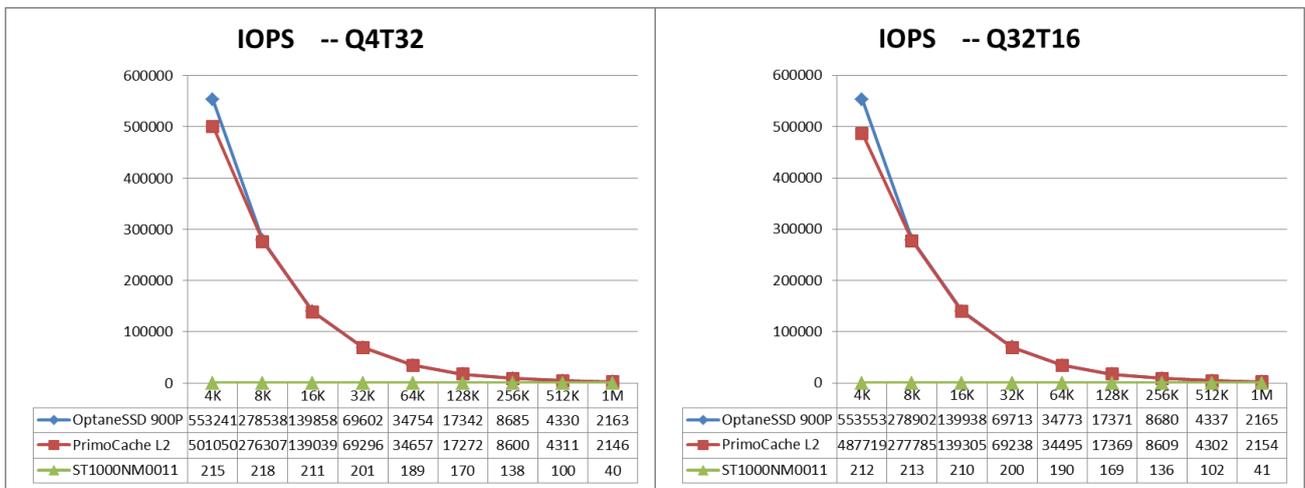
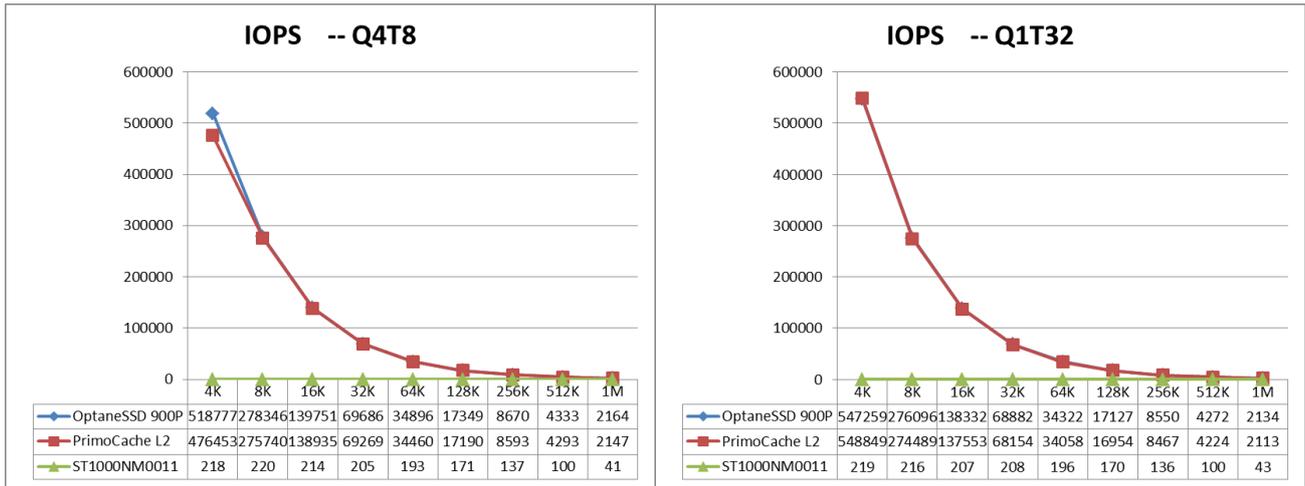
从单线程单队深的延迟测量结果可以看到 SSD 900P 4KB 随机读的最低平均延迟为 **11 μs**，基本和官方标称 10 μs 一致。PrimoCache 缓存的平均延迟为 13 μs，比 SSD 900P 多了 2 μs 的额外处理开销。



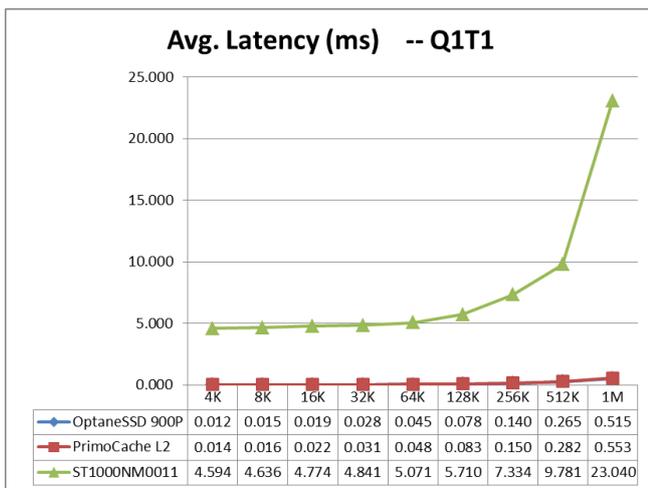
## 随机标准写

随机标准写的测试结果和随机读类似，SSD 900P 在单线程单队深下的 4KB 随机写 IOPS 数值接近 8 万，随着线程数或队深数的增加，其 IOPS 数值显著上升，最高可达到 55 万以上，高于官方标称的 50 万。





从单线程单队深的延迟测量结果上可以看到 SSD 900P 4KB 随机写的最低平均延迟为 **12 μs**，略高于官方标称的 10 μs。PrimoCache 缓存的平均延迟为 14 μs，和 4KB 随机读一样，也是比 SSD 900P 多了 2 μs 的额外处理开销。



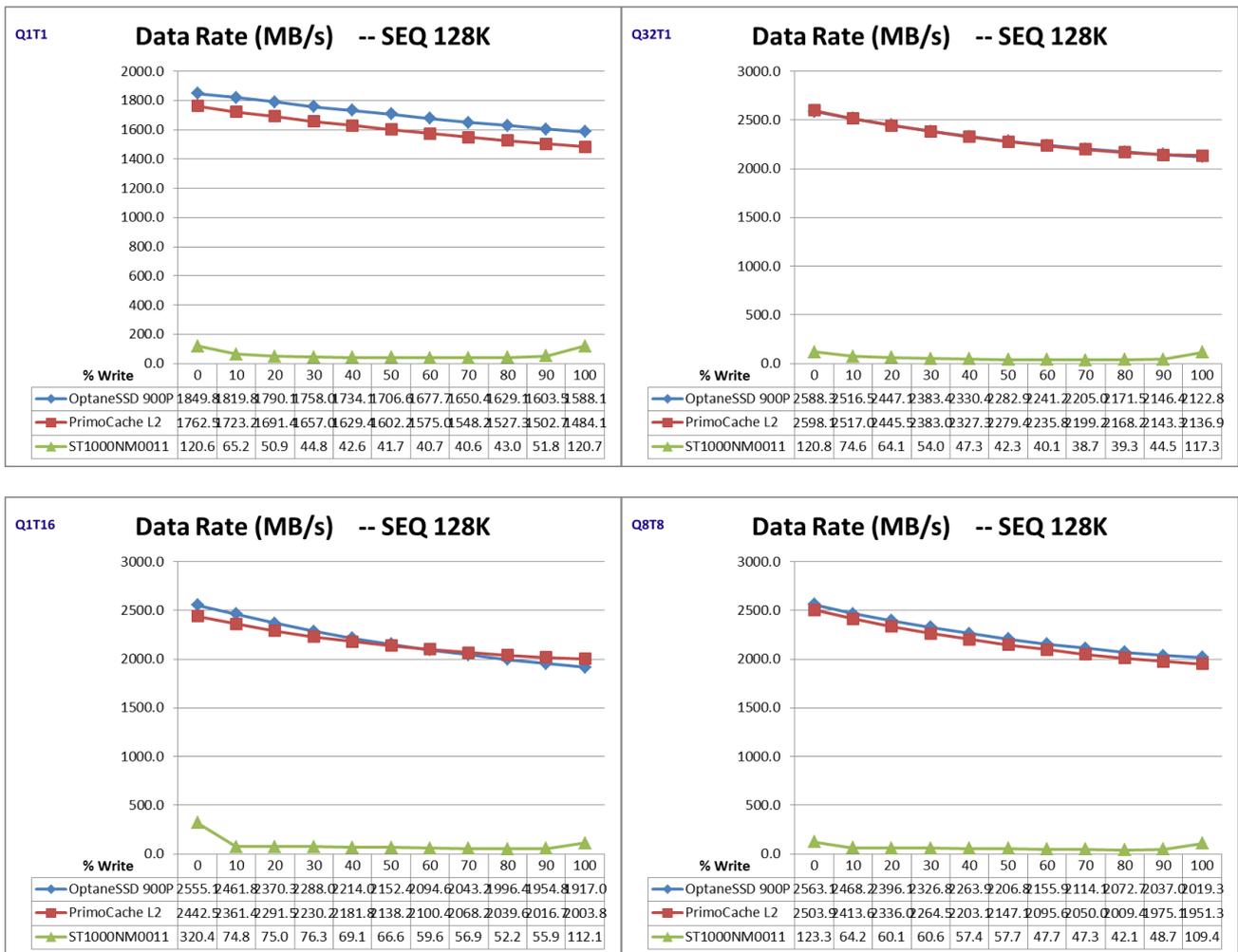
# 混合读写测试

混合读写测试主要是测试存储设备同时进行读写操作的性能。本文测试了不同数据块大小在不同线程数不同队深下的顺序混合读写和随机混合读写的性能。这里选取部分测试结果图表进行说明。图表中的横坐标表示“%写”，即测试中写 IO 数目占全部 IO 数目的百分比。比如“30%写”表示测试 IO 由 30%写 IO 和 70%读 IO 组成。“0%写”是完全读（即标准读），“100%写”则是完全写（即标准写）。

## 顺序混合读写

以下图表是 128KB 数据块大小分别在 Q1T1、Q32T1、Q1T16 和 Q8T8 下顺序混合读写的测试结果。可以看到随着读写比例的变化，SSD 900P 对于 128KB 顺序混合读写性能基本按照完全读写性能的比例呈现相应的线性变化。这是比较惊人的测试结果，很多其它 SSD 或机械硬盘都存在完全读写性能较高但混合读写性能急剧下降的现象。

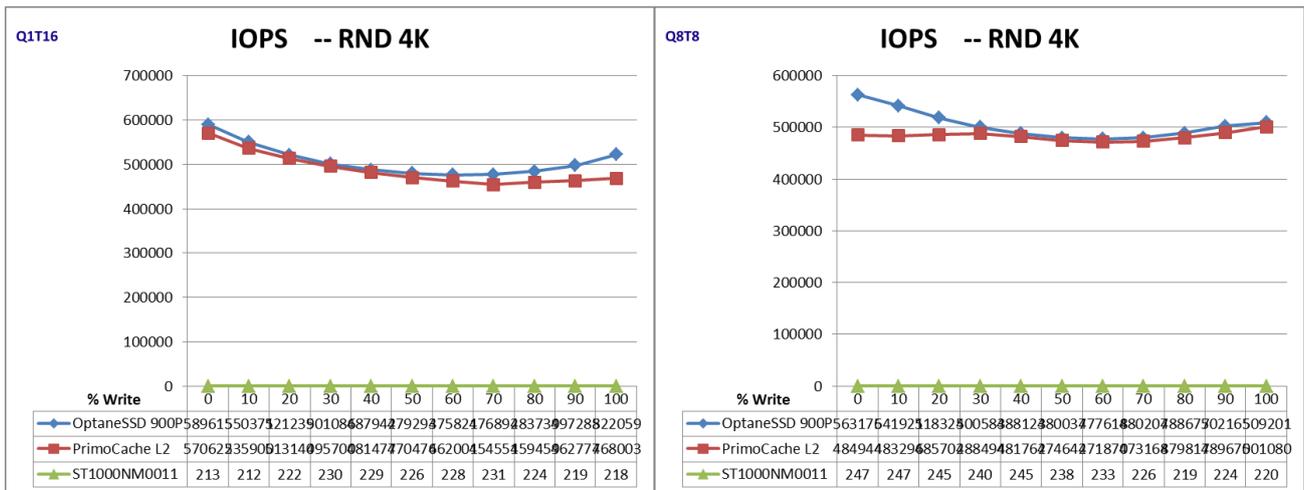
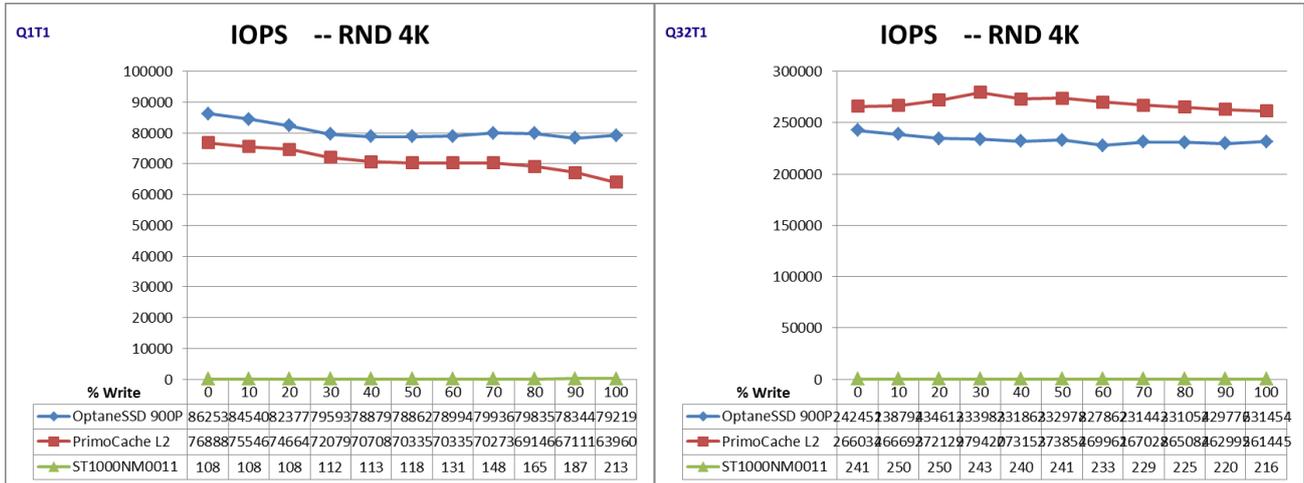
和标准读写测试结果类似，PrimoCache 缓存在顺序混合读写上的性能也是基本接近或达到 Optane SSD 900P 的性能。



## 随机混合读写

以下图表显示了 4KB 数据块大小分别在 Q1T1、Q32T1、Q1T16 和 Q8T8 下随机混合读写的测试结果。可以看到随着读写比例的变化，SSD 900P 在低线程数时还是基本呈现线性变化的，但在高线程数时表现出曲线变化，混合读写的 IOPS 值低于其相应的线性值，不过偏离并不大。

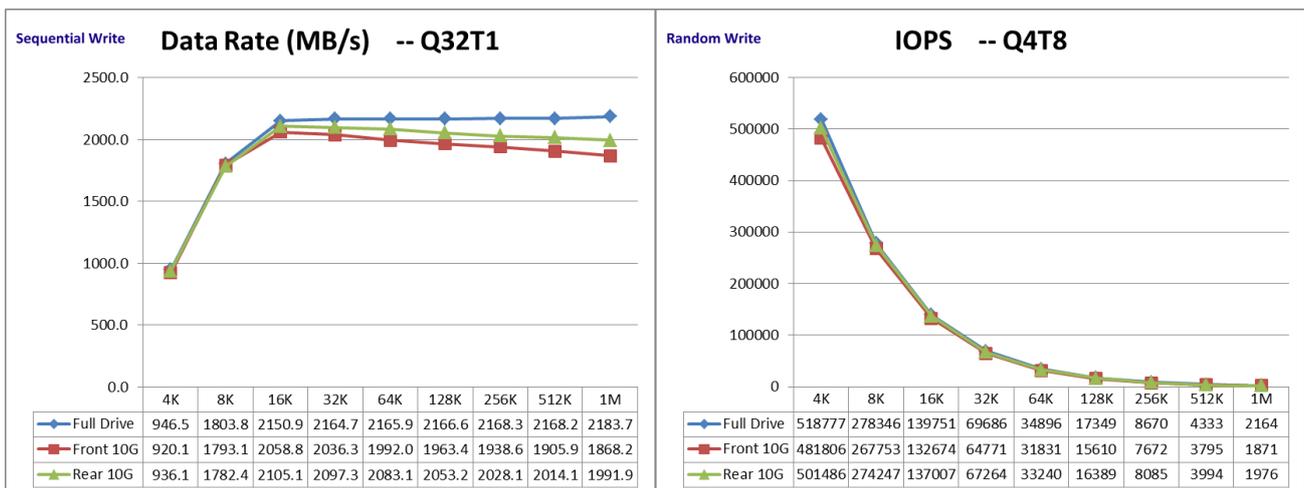
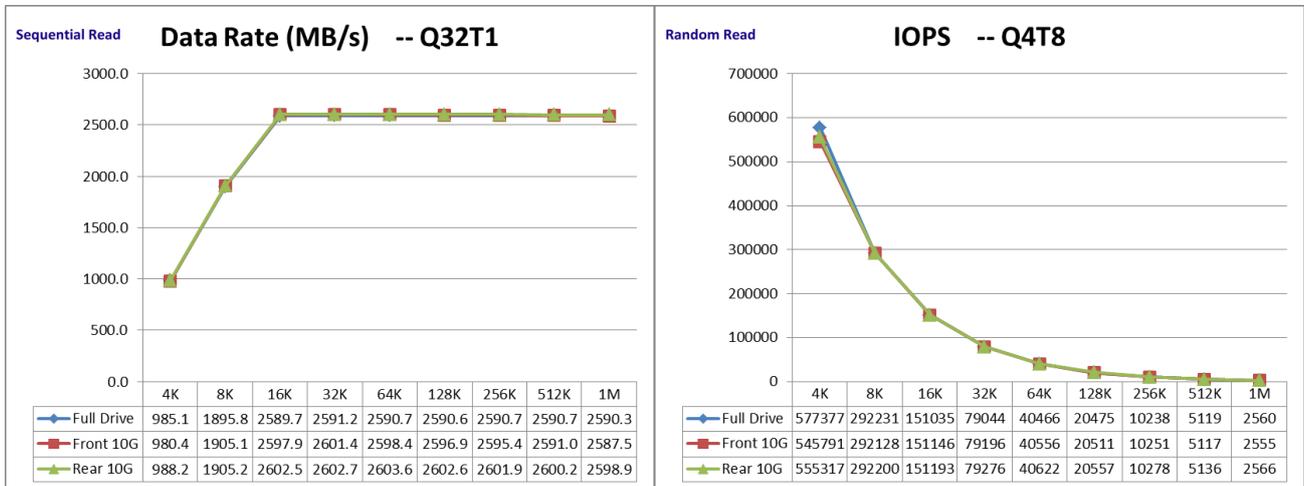
PrimoCache 缓存在 4KB 随机混合读写上的表现也和 4KB 标准读写的测试结果类似，这里不再赘述。



## 全盘和区间测试

为评估存储设备在不同物理地址空间上的性能差异，本文对 Optane SSD 900P 的全部地址空间和前后端 10GB 范围的地址空间分别进行了标准性能测试并进行比较。所有测试均是在文件填满整个 SSD 的条件下完成。测试结果对比图表中的“Full Drive”指全盘空间，“Front 10G”指前端 10GB 空间，实际测试选取了物理地址 13GB-23GB 的空间，“Rear 10G”指后端 10GB 空间，实际测试选取的是物理地址 223GB-233GB 的空间。

不同线程不同队深下的测试结果基本类似，因此这里仅选取一些典型图表进行分析。从这些图表可以看到，顺序读性能在整个物理地址空间具有非常好的一致性，顺序写性能在不同物理地址空间上略有些差异，4KB 随机读写 IOPS 也随着地址空间不同略微存在差异，但 8KB 或其它大小的数据块随机读写性能基本一致。整体上来说，其空间一致性是相当好的。



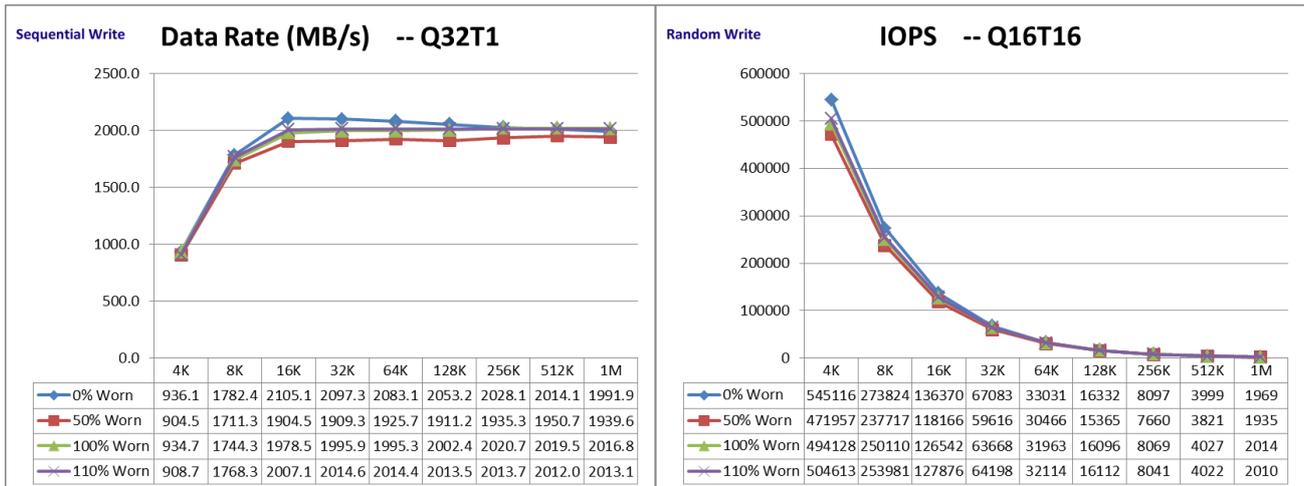
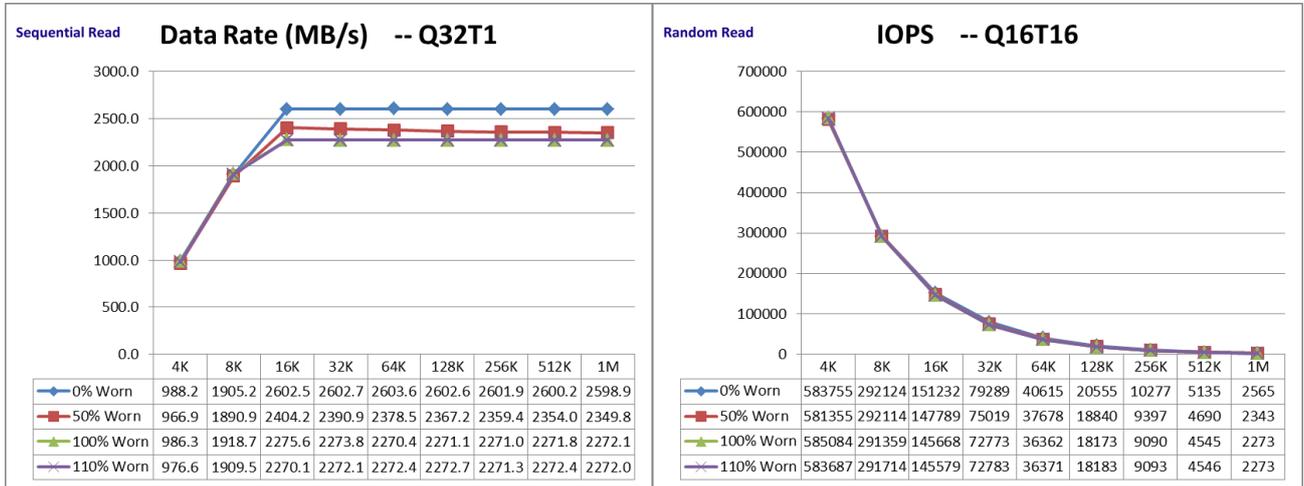
## 寿命影响测试

SSD 通常有写入次数的限制，寿命影响测试主要评估 SSD 写入寿命对性能的影响，即测试大量多次写入后的性能变化。由于对全盘进行寿命测试非常耗时，即使以最大速度不间断写入达到 Optane SSD 900P (280GB) 的标称寿命，所需时间也超过 30 天，因此实际测试仅选取其中 10GB 地址空间进行。SSD 900P 的标称写入寿命为 18.69TB 每 GB 容量，因此 10GB 容量的写入寿命为 186.9TB。在此空间均匀写入 0TB (0%寿命)、90TB (50%寿命)、180TB (100%寿命)、200TB (110%寿命) 后分别进行标准性能测试。

从测试结果对比图表中可以发现，顺序读性能（以 Q32T1 为例）在 50%寿命后从 2600MB/s 减少到 2350MB/s，仅下降不到 10%，在达到 100%寿命或更多时，性能基本稳定在 2270MB/s，仅下降不到 13%。4KB 随机读 IOPS 性能（以 Q16T16 为例）即使在达到 110%寿命也仍然保持不变，令人惊

讶。顺序写性能（Q32T1）和 4KB 随机写 IOPS（Q16T16）在 50%寿命后分别最高仅有大约 10%和 13%的性能降低。比较奇怪的是在达到 100%寿命后，写入性能反而比 50%寿命时还要好一些。

在达到 110%寿命后，我们对这块地址空间进行了读写数据准确性验证。经过多次测试，写入后并断电回读的数据与源数据一致，并无出现错误。因此从测试结果来看，Intel 标称的寿命值还是比较保守的，实际寿命要超过标称值。



## 结论

综合上述测试结果，可以看到，Intel Optane SSD 900P 各项性能指标的实测结果都能达到官方标称值，甚至大部分结果都超出官方值。其在读写速度、IOPS、访问延迟上的性能都非常强劲，远远超出当前市场主流 SSD。尤其在中高工作负载下，其性能可以达到充分利用。

Optane SSD 900P 在混合读写测试、不同物理地址空间测试以及在不同使用寿命阶段都保持了非常好的性能一致性，没有显著的性能降低，即使在达到标称写入寿命后，所测性能也仅下降 10%左右，不超过 15%。

在使用 PrimoCache 软件将 SSD 900P 用作其它硬盘的缓存后，再对被缓存的硬盘进行测试可以看到，其性能表现基本接近或达到 SSD 900P 的性能。在某些情形下，甚至还超出 SSD 900P 的自身性能。

结合 SSD 900P 的优异读写性能及超高稳定性和 PrimoCache 软件的缓存性能，如前所述，除了桌面消费级市场上的应用，SSD 900P 也适合用作工作站或服务器上大容量低速硬盘的缓存，简单方便即可提升这些硬盘的读写速度和 IO 性能，使其接近或达到，甚至超过 SSD 900P 的性能。非常适合有预算控制或不希望修改系统配置迁移数据的用户。

本文虽然没有对 PrimoCache 缓存在服务器上具体应用的性能提升进行测试，但本文的测试方案基本涵盖了存储设备在具体应用中可能遇到的各种 IO 处理场景，体现了存储设备的基准性能。以 SQL Server 为例，下表列出了各种可能的 SQL Server 操作。因此要评估 SQL Server 系统上的磁盘 IO 性能，只要测试其相应数据块大小的随机/顺序读写性能即可得到一个定量的结果。

文件类型	操作	读数据块	写数据块	线程数	I/O 类型
数据文件	日常活动	8KB - 128KB	8KB - 128KB	基于 MaxDOP	随机
	Checkpoint	N/A	64KB - 128KB	CPU 插座数	随机
	LazyWriter	N/A	64KB - 128KB	每 NUMA 节点 1 个	随机
	Bulk Insert	N/A	8KB - 128KB	基于 MaxDOP	顺序
	备份	1 MB	1 MB	基于 MaxDOP	顺序
	恢复	64KB	64KB	基于 MaxDOP	顺序
	DBCC Checkdb (不带修复)	8KB - 64KB	N/A	基于 MaxDOP	顺序
	Rebuild Index	最高至 512 KB	8KB - 128KB	基于 MaxDOP	顺序
	ReadAhead	最高至 512 KB	N/A	基于 MaxDOP	顺序
日志文件	日常活动	512B - 64KB	512B - 64KB	每 NUMA 软节点 1 个日志写入线程，最多 4 个	顺序

(注：此表引自 Microsoft Diskspd 帮助文档)

值得一提的是，PrimoCache 也支持同时将内存作为硬盘的缓存，为了准确反映使用 Optane SSD 900P 作为缓存的性能，本文测试 PrimoCache 缓存性能时没有增加设置内存作为缓存，而是纯粹使用 SSD 900P。在实际应用中，可以增加设置一部分内存作为一级缓存，SSD 900P 作为二级缓存，以获取更佳性能。

当然，为了测试 PrimoCache 缓存的最大性能，本文所做测试都是在保证测试文件完全被缓存至 SSD 900P 的前提下进行。在实际应用中，由于缓存空间相对于被缓存盘的容量通常会比较有限，不能保证所有文件都会被缓存，此外根据实际需求，缓存配置也有所不同，因此实际应用中的性能可能和本文测试结果会有差异。

## 应用场景

从上述结论中不难发现，使用 PrimoCache 软件将 Optane SSD 900P 作为系统低速存储设备的高速大容量缓存，无论是系统的读取性能还是写入性能都可以得到大幅提升，因此可应用场景和领域是非常宽泛的。对于当前存在读取性能或写入性能瓶颈的工作站/中小型服务器系统，Optane SSD 900P + PrimoCache 即可轻松解决读写瓶颈问题。如果要新搭建一个读取或写入性能有高要求且存储数据量大的系统，那么 SSD 900P 缓存+大容量低速硬盘的方案相对于直接使用大容量高速硬盘的方案具有超高性价比。下面简单列举了一些 SSD 900P 缓存方案可应用的工作场景，但并不局限于这些场景。

- 工作站：例如对于一台主要用于 3DMAX、MAYA、C4D 等 3D 软件所制作动画进行渲染输出的图形工作站，应用 SSD 900P 缓存方案后，可以大幅度提升整个系统的 IO 效率，大大减少渲染时间。
- 针对渲染农场的应用：在渲染过程中，渲染节点会对服务器产生大量的，近乎同时的 I/O 请求。尤其是在节点比较多的情况下，常规存储系统的 I/O 性能很难满足大规模渲染应用的需求。而应用 SSD 900P 缓存方案则可以很好地解决这个问题。
- 针对 VDI 的应用：大批瘦客户机或无盘站同时启动时，将对服务器端的磁盘 IO 产生大量的读操作，而 SSD 900P 高带宽、低延时、高寿命等特点，结合 PrimoCache 软件，可以提升整个 VDI 解决方案的性价比。
- 传统 NAS 或分布式存储方案：采用 SSD 900P 缓存方案加速存储性能，可以简单安全的实现系统升级。
- 网吧服务器：例如当前流行的绝地求生游戏（俗称“吃鸡”游戏）会产生大量的数据回写，传统的 SSD 由于寿命和容量的限制，很快就会达到 SSD 极限。而 SSD 900P 超高寿命超快写入的特点，作为网吧服务器的回写盘缓存，有效解决网吧“吃鸡”的烦恼！

## 附件：测试结果报告

### 标准读写测试结果报告

[顺序标准读测试](#)

[顺序标准写测试](#)

[随机标准读测试](#)

[随机标准写测试](#)

### 混合读写测试结果报告

[Q1T1 测试](#)

[Q32T1 测试](#)

[Q1T16 测试](#)

[Q8T8 测试](#)

### 全盘和区间测试结果报告

[顺序标准读测试](#)

[顺序标准写测试](#)

[随机标准读测试](#)

[随机标准写测试](#)

### 寿命影响测试结果报告

[顺序标准读测试](#)

[顺序标准写测试](#)

[随机标准读测试](#)

[随机标准写测试](#)