



BUILDING AN NVIDIA DIGITS DEVBOX

DG-07680-001_v01 | May 2015

Design Guide



DOCUMENT CHANGE HISTORY

DG-07680-001_v01

Version	Date	Authors	Description of Change
01	May 29, 2015	CP, TS	Initial release

TABLE OF CONTENTS

- Building an NVIDIA DIGITS DevBox..... 5**
 - Design Requirements 6
 - Hardware..... 6
 - Motherboard..... 6
 - GPUs 7
 - CPU and Cooling 8
 - Memory and Storage..... 8
 - Chassis, Thermal, and Acoustic Considerations..... 9
 - Power Supply12
 - Operating and Software Installation12

LIST OF FIGURES

Figure 1.	Topology Used in the Sample NVIDIA DIGITS DevBox.....	7
Figure 2.	Corsair Carbide Series Air 540 Chassis (with upgraded front fans).....	9
Figure 3.	Custom Bracket to Support GPUs	10
Figure 4.	Custom GeForce GTX Titan X Bracket.....	11
Figure 6.	EVGA SuperNOVA 1600 W P2 Power Supply	12

LIST OF TABLES

Table 1.	GPU Comparisons.....	7
Table 2.	Recommended CPUs	8

BUILDING AN NVIDIA DIGITS DEVBOX

Applied research in deep learning requires the fastest possible experiment turnaround times to rapidly explore multiple network architectures and manipulate and curate datasets to reduce solution delivery times for internal and external customers. Effective, deep learning training requires a system that is designed and built for this specific task. Simply adding RAM and GPUs to an ordinary workstation can cause serious issues with stability, cooling, PCIe topology, and total system power. The NVIDIA® DIGITS™ DevBox deep learning system was designed and built to the high-performance standards to meet these needs.

This document describes the key considerations and decisions that we made at each step of the process to help guide you through your own design & build. It includes the following general areas:

- ▶ Design requirements
- ▶ Hardware
- ▶ Operating system
- ▶ Software



Note: NVIDIA is sharing our experience to give you some suggestions about how to go about building your own NVIDIA DIGITS DevBox. This document is not intended as step-by-step instructions. We are therefore unable to provide technical support or service to the DIY community based on the information contained herein.

DESIGN REQUIREMENTS

Our goal was to build the fastest machine learning training device that can be deployed in a standard office or research environment. This key requirement forced us to make some fundamental decisions about maximum power consumption while also providing adequate cooling without sacrificing performance and without generating excessive noise. We settled on the following general objectives:

- ▶ Highest deep learning throughput
- ▶ 1600W peak power limit (limit of a standard household 15A circuit)
- ▶ Sufficient cooling to protect electronics and dissipate system heat
- ▶ Reasonably quiet acoustic performance even under peak load

HARDWARE

The choice of hardware greatly impacts system stability and performance; especially when relying on consumer-grade products for cost-effectiveness.

Motherboard

Effective deep learning requires multiple GPUs. However, suitable PCIe topology is critical to being able to use those GPUs efficiently. Synchronous Stochastic Gradient Descent (SGD) for deep learning relies on broadcast communication between the GPUs. SGD acceleration needs P2P DMAs to work between devices. This means that all GPUs *must* be on the same I/O hub with a very fast PCIe switches. Workstation motherboards based on the Intel X99 chipset with a PLX bridge setup can support four PCIe Generation3 x16 cards at either full speed or with minimal drop-off. Performance depends on communications between GPU cards and overall traffic patterns.



Note: Our sample build used the *ASUS X99-E WS workstation motherboard* that supports Intel LGA 2011-v2 CPUs while drawing only 20W.



GPUs

NVIDIA Maxwell GPUs deliver 20% more GFLOPS/watt than Kepler-based GPUs. Of these, the NVIDIA GM200 GPU offers both the fastest single-precision performance (the critical metric for training DNNs and other deep learning tasks) and best power efficiency. The GM200 is available in the GeForce™ Titan X with 12GB of memory per board, which allows the system to handle larger models. Table 1 lists the GPUs.

Our sample build used four NVIDIA GeForce GTX Titan X cards that consume a combined of 1,000 W; leaving 580 W for the remaining system components.

Figure 1 shows the system topology used.

Table 1. GPU Comparisons

GPU	Code Name	Family	Transistors	Single Precision (GFLOPs)	TDP (W)	GFLOPS/W
GeForce GTX 780 Ti	GK110	Kepler	7B	5046	250	20.2
GeForce GTX Titan Black	GK110	Kepler	7B	5121	250	20.5
GeForce GTX 980	GM204	Maxwell	5B	4612	165	28.0
GeForce GTX Titan X	GM200	Maxwell	8B	6605	250	26.4

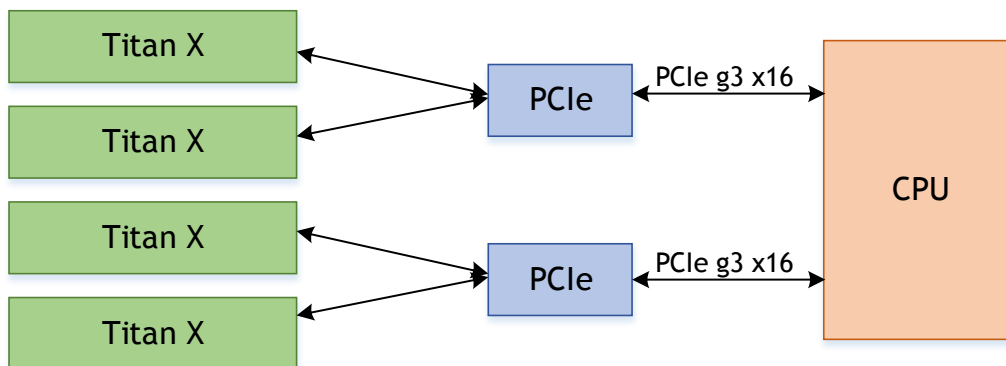


Figure 1. Topology Used in the Sample NVIDIA DIGITS DevBox

CPU and Cooling

The CPU must be compatible with the selected chipset while providing sufficient PCIe support. Consumer CPUs are ideal because the application target is highly fault tolerant and the NVIDIA DIGITS DevBox is acting as a workstation instead of a server. Table 2 lists the recommended CPUs.

Table 2. Recommended CPUs

CPU	PCIe Configuration	Cores	TDP (W)
Core i7-5930K	40	6	140W
Core i7-5960K	40	8	140W

Pairing the CPU with effective cooling is crucial for optimal performance, especially when the GPUs are under peak load. For the NVIDIA DIGITS DevBox, we chose the *Intel Core i7-5930K CPU* with a Corsair Hydro H60 cooler.

Memory and Storage

RAM is important for handling large DNN files and datasets. The Intel Core i7-5930K CPU can stably support up to 64 GB RAM. An Intel Xeon processor can handle more RAM and allow ECC. However, using the Intel Xeon processor will significantly increase the cost.

The mass storage I/O subsystem is also important because the GPUs can consume data at rates far beyond the capabilities of a single mechanical drive. SSDs provide speed but increase system cost while limiting capacity. We decided to leverage the best of both worlds by combining multiple mechanical drives into a RAID array and using SSDs for cache and boot. This combination delivers optimal performance while containing costs.

Our sample build uses three 3-TB disks combined in a RAID5 array (~250 MB/s sustained read) and a single 512 MB m.2 SSD-based cache (1GB/s) for a total of 6 TB of usable storage with very fast caching. A separate 256GB SSD serves as the system boot drive. The system includes an Icy-Dock FlexCage in the two 5.25-inch bays for easy hot-swapping of the RAID drives.

Chassis, Thermal, and Acoustic Considerations

Acoustics and heat management are major considerations, especially when deploying the NVIDIA DIGITS DevBox in a normal office environment. A chassis that separates the power supply and disks from the heat generated by the CPU and GPUs is ideal. Further, effective cooling requires generating significant positive air pressure inside the chassis to facilitate drawing air through the gaps between the closely-spaced NVIDIA GeForce GTX Titan X cards. These cards include blower-style cooling fans that direct airflow through the GPU heat sinks and out of the chassis to prevent internal heat buildup.

Our sample build used a *Corsair Carbide Series Air 540* chassis with upgraded front fans. We also chose to close off unneeded vents in order to increase air pressure within the chassis.



Figure 2. Corsair Carbide Series Air 540 Chassis (with upgraded front fans)

Spacing between the GPU cards is critical for consistent airflow. This isn't usually a problem in a chassis where GPUs are mounted vertically. However, the horizontally-suspended GPU cards in the Corsair chassis can sag against the lower cards and impinge airflow. To solve this problem, we built a custom rear support bracket to ensure proper GPU spacing and recommend that you consider doing the same. The bracket we built is shown in Figure 3.



Figure 3. Custom Bracket to Support GPUs

The DevBox also uses a special version of the GeForce GTX Titan X that includes a modified bracket which removes unused display I/O ports to maximize airflow (Figure 4).



Original Bracket



Custom Bracket

Figure 4. Custom GeForce GTX Titan X Bracket



CAUTION: MODIFYING THE NVIDIA GRAPHICS CARDS WILL VOID YOUR WARRANTY.

Fan tuning software can help maximize airflow while keeping the system within acoustic guidelines.

We expect that using off-the-shelf components that do not include the thermal and acoustic optimizations described will cause a 15% drop in performance.

Power Supply

The power supply should provide enough power to operate the system components along with some headroom to ensure stable operation. The total dissipated power for all of the system components we used in our sample build is between 1,200 and 1,300 watts.

Our sample build uses an EVGA SuperNOVA 1600W P2 power supply that delivers approximately 90% efficiency at 100% load (1,400 watts), ensuring system stability at peak workloads.



Figure 5. EVGA SuperNOVA 1600 W P2 Power Supply

OPERATING AND SOFTWARE INSTALLATION

Following is a list of the operating system and software used in our NVIDIA DIGITS DevBox build:

- ▶ Ubuntu 14.04 operating system available from <http://releases.ubuntu.com/14.04/>
- ▶ cuDNN library available from <https://developer.nvidia.com/cuda-registered-developer-program> (registration required)

The NVIDIA DIGITS DevBox ships with several deep learning software packages that can be built from the following sources:

- ▶ NVIDIA DIGITS is an interactive environment for training, evaluating, and experimenting with neural networks. The source for version 1.0.3 installed in the sample NVIDIA DIGITS DEVBOX described in this article is available from <https://github.com/NVIDIA/DIGITS/releases/tag/v1.0.3>
- ▶ The installed version of Caffe is from an NVIDIA branch available from <https://github.com/NVIDIA/caffe/releases/tag/v0.10.0>
- ▶ The Torch version is from <https://github.com/torch/distro> using a specific commit: <https://github.com/torch/distro/archive/e0c565120622f99ef6e1ca7fcca66cfe2da34fc>

- ▶ The packaged Theano 0.7.0 version is available from <https://pypi.python.org/pypi/Theano/0.7.0>
- ▶ The BIDMach 1.0.0 release tarball is available from <http://bid2.berkeley.edu/bid-data-project/download/>



Note: NVIDIA has only tested this software on our DIGITS DevBox build and cannot guarantee that this will work on your system. Each listing includes a link for downloading the source and/or binary files.

Notice

The information provided in this specification is believed to be accurate and reliable as of the date provided. However, NVIDIA Corporation ("NVIDIA") does not give any representations or warranties, expressed or implied, as to the accuracy or completeness of such information. NVIDIA shall have no liability for the consequences or use of such information or for any infringement of patents or other rights of third parties that may result from its use. This publication supersedes and replaces all other specifications for the product that may have been previously supplied.

NVIDIA reserves the right to make corrections, modifications, enhancements, improvements, and other changes to this specification, at any time and/or to discontinue any product or service without notice. Customer should obtain the latest relevant specification before placing orders and should verify that such information is current and complete.

NVIDIA products are sold subject to the NVIDIA standard terms and conditions of sale supplied at the time of order acknowledgement, unless otherwise agreed in an individual sales agreement signed by authorized representatives of NVIDIA and customer. NVIDIA hereby expressly objects to applying any customer general terms and conditions with regard to the purchase of the NVIDIA product referenced in this specification.

NVIDIA products are not designed, authorized or warranted to be suitable for use in medical, military, aircraft, space or life support equipment, nor in applications where failure or malfunction of the NVIDIA product can reasonably be expected to result in personal injury, death or property or environmental damage. NVIDIA accepts no liability for inclusion and/or use of NVIDIA products in such equipment or applications and therefore such inclusion and/or use is at customer's own risk.

NVIDIA makes no representation or warranty that products based on these specifications will be suitable for any specified use without further testing or modification. Testing of all parameters of each product is not necessarily performed by NVIDIA. It is customer's sole responsibility to ensure the product is suitable and fit for the application planned by customer and to do the necessary testing for the application in order to avoid a default of the application or the product. Weaknesses in customer's product designs may affect the quality and reliability of the NVIDIA product and may result in additional or different conditions and/or requirements beyond those contained in this specification. NVIDIA does not accept any liability related to any default, damage, costs or problem which may be based on or attributable to: (i) the use of the NVIDIA product in any manner that is contrary to this specification, or (ii) customer product designs.

No license, either expressed or implied, is granted under any NVIDIA patent right, copyright, or other NVIDIA intellectual property right under this specification. Information published by NVIDIA regarding third-party products or services does not constitute a license from NVIDIA to use such products or services or a warranty or endorsement thereof. Use of such information may require a license from a third party under the patents or other intellectual property rights of the third party, or a license from NVIDIA under the patents or other intellectual property rights of NVIDIA. Reproduction of information in this specification is permissible only if reproduction is approved by NVIDIA in writing, is reproduced without alteration, and is accompanied by all associated conditions, limitations, and notices.

ALL NVIDIA DESIGN SPECIFICATIONS, REFERENCE BOARDS, FILES, DRAWINGS, DIAGNOSTICS, LISTS, AND OTHER DOCUMENTS (TOGETHER AND SEPARATELY, "MATERIALS") ARE BEING PROVIDED "AS IS." NVIDIA MAKES NO WARRANTIES, EXPRESSED, IMPLIED, STATUTORY, OR OTHERWISE WITH RESPECT TO THE MATERIALS, AND EXPRESSLY DISCLAIMS ALL IMPLIED WARRANTIES OF NONINFRINGEMENT, MERCHANTABILITY, AND FITNESS FOR A PARTICULAR PURPOSE. Notwithstanding any damages that customer might incur for any reason whatsoever, NVIDIA's aggregate and cumulative liability towards customer for the products described herein shall be limited in accordance with the NVIDIA terms and conditions of sale for the product.

VESA DisplayPort

DisplayPort and DisplayPort Compliance Logo, DisplayPort Compliance Logo for Dual-mode Sources, and DisplayPort Compliance Logo for Active Cables are trademarks owned by the Video Electronics Standards Association in the United States and other countries.

HDMI

HDMI, the HDMI logo, and High-Definition Multimedia Interface are trademarks or registered trademarks of HDMI Licensing LLC.

OpenCL

OpenCL is a trademark of Apple Inc. used under license to the Khronos Group Inc.

Trademarks

NVIDIA, the NVIDIA logo, GeForce, and DIGITS are trademarks and/or registered trademarks of NVIDIA Corporation in the U.S. and other countries. Other company and product names may be trademarks of the respective companies with which they are associated.

Copyright

© 2015 NVIDIA Corporation. All rights reserved.