



# **INTEL DISTRIBUTION FOR PYTHON – ADVANTAGES AND ACCELERATION OF MACHINE LEARNING WORKLOADS**

Intel Software Developer Conference – London, 2017

# Legal Disclaimer and Optimization Notice

INFORMATION IN THIS DOCUMENT IS PROVIDED “AS IS”. NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. INTEL ASSUMES NO LIABILITY WHATSOEVER AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO THIS INFORMATION INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT.

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products.

OpenCL and the OpenCL logo are trademarks of Apple Inc. used by permission by Khronos.

Copyright © 2017, Intel Corporation. All rights reserved. Intel, Pentium, Xeon, Xeon Phi, Core, VTune, Cilk, and the Intel logo are trademarks of Intel Corporation in the U.S. and other countries.

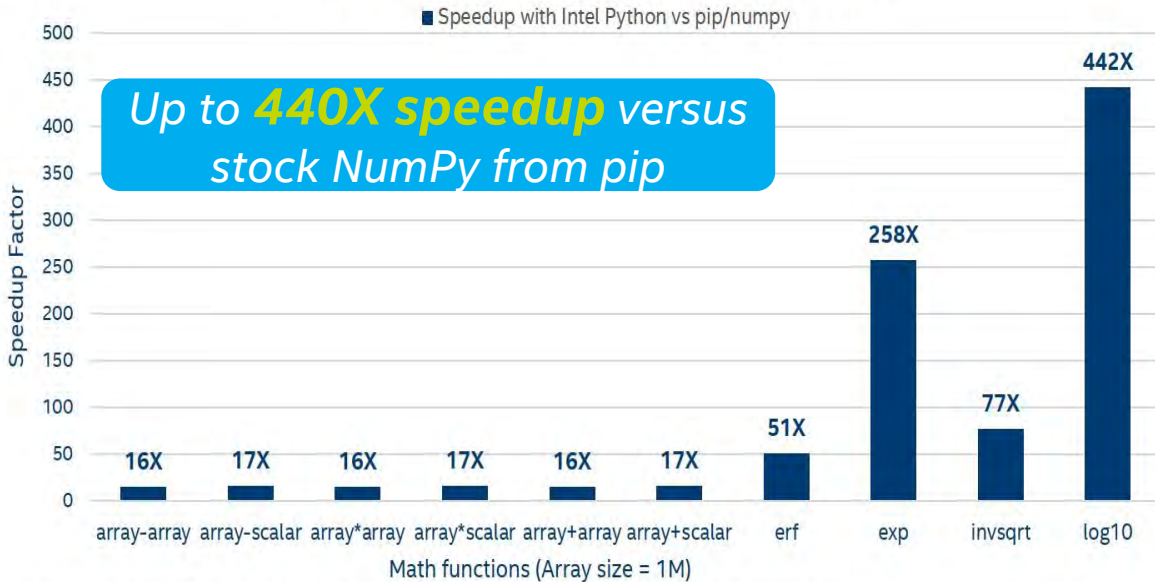
## Optimization Notice

Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.

Notice revision #20110804

# Faster Python\* with Intel® Distribution for Python 2018

Intel® Distribution for Python\* Performance Speedups  
for Select Math Functions on Intel® Xeon™ Processors



**Configuration:** Hardware: Intel® Xeon® CPU E5-2699 v4 @ 2.20GHz (2 sockets, 22 cores per socket, 1 thread per core – HT is off), 256GB DDR4 @ 2400MHz.  
Software: Stock CentOS Linux release 7.3.1611 (Core), python 3.6.2, pip 9.0.1, numpy 1.13.1, scipy 0.19.1, scikit-learn 0.19.0. Intel® Distribution for Python\* 2018 Gold: mkl 2018.0.0 Intel\_4, daal 2018.0.0 20170814, numpy 1.13.1 py36\_intel\_15, openmp 2018.0.0 intel\_7, scipy 0.19.1 np113py36\_intel\_11, scikit-learn 0.18.2 np113py36\_intel\_3

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more information go to <https://www.intel.com/performance>. Benchmark Source: Intel Corporation.

**Optimization Notice:** Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice. Notice revision #20110804. \*Other brands and names are the property of their respective owners.

Learn More: [software.intel.com/distribution-for-python](https://software.intel.com/distribution-for-python)

## High Performance Python Distribution

- Accelerated NumPy, SciPy, scikit-learn well suited for scientific computing, machine learning & data analytics
- Drop-in replacement for existing Python. No code changes required
- Highly optimized for latest Intel processors
- Take advantage of [Priority Support](#) – connect direct to Intel engineers for technical questions<sup>2</sup>

## What's New in 2018 version

- Updated to latest version of Python 3.6
- Optimized scikit-learn for machine learning speedups
- Conda build recipes for custom infrastructure

Software & workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark & MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information & performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more information go to <http://www.intel.com/performance>.

### Optimization Notice

<sup>2</sup>Paid versions only.

Copyright © 2017, Intel Corporation. All rights reserved.

\*Other names and brands may be claimed as the property of others.

Optimization Notice: Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.





# Python\* Landscape

## Challenge#1

Domain experts are not professional software programmers

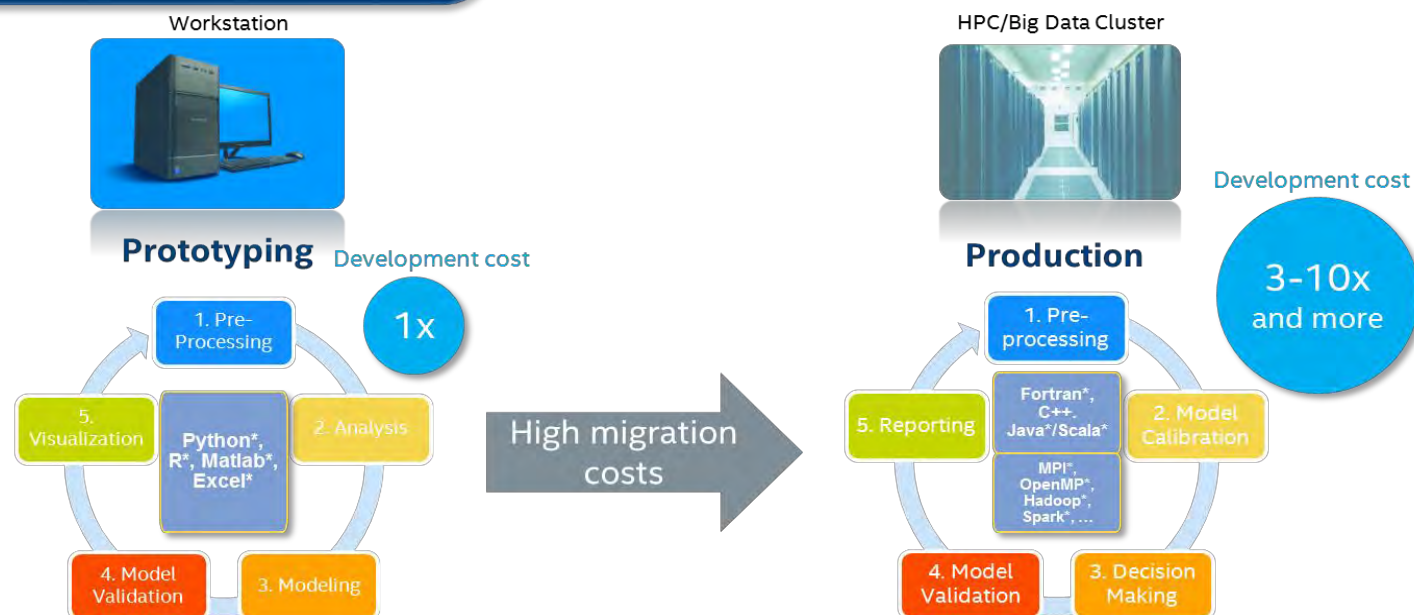
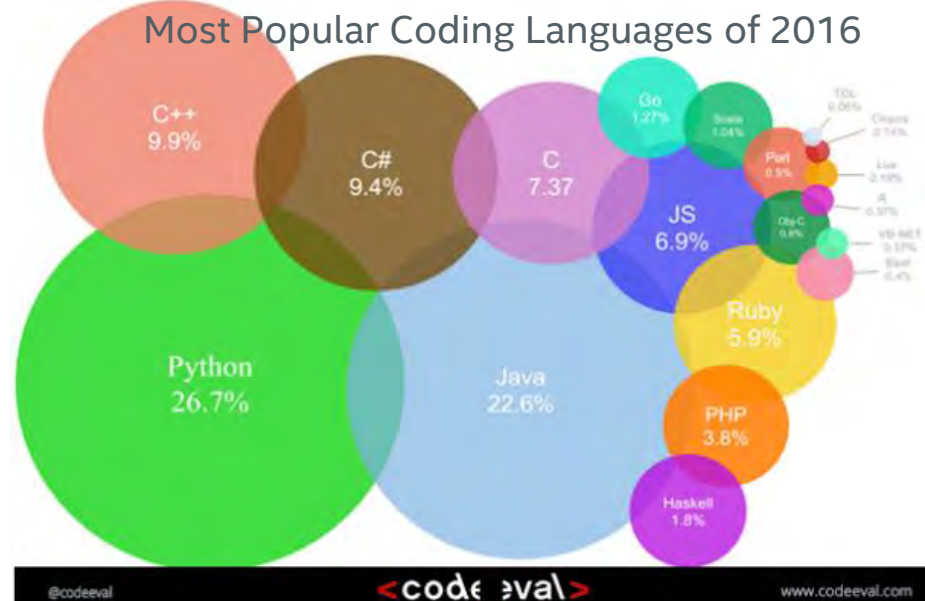
## Challenge#2

Python performance limits migration to production systems

### Intel's Python Tools

- Accelerate Python performance
- Enable easy access
- Empower the community

**Adoption of Python**  
continues to grow among domain experts & developers for its productivity benefits



### Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.

\*Other names and brands may be claimed as the property of others.

# What's Inside Intel® Distribution for Python

High Performance Python\* for Scientific Computing, Data Analytics, Machine Learning

FASTER PERFORMANCE	GREATER PRODUCTIVITY	ECOSYSTEM COMPATIBILITY
<b>Performance Libraries, Parallelism, Multithreading, Language Extensions</b>	<b>Prebuilt &amp; Accelerated Packages</b>	<b>Supports Python 2.7 &amp; 3.6, conda, pip</b>
<p>Accelerated NumPy/SciPy/scikit-learn with Intel® MKL<sup>1</sup> &amp; Intel® DAAL<sup>2</sup></p> <p>Data analytics, machine learning &amp; deep learning with scikit-learn, pyDAAL</p> <p>Scale with Numba* &amp; Cython*</p> <p>Includes optimized mpi4py, works with Dask* &amp; PySpark*</p> <p>Optimized for latest Intel® architecture</p>	<p>Prebuilt &amp; optimized packages for numerical computing, machine/deep learning, HPC, &amp; data analytics</p> <p>Drop in replacement for existing Python - No code changes required</p> <p>Jupyter* notebooks, Matplotlib included</p> <p>Conda build recipes included in packages</p> <p>Free download &amp; free for all uses including commercial deployment</p>	<p>Compatible &amp; powered by Anaconda*, supports conda &amp; pip</p> <p>Distribution &amp; individual optimized packages also available at conda &amp; Anaconda.org, YUM/APT, Docker image on DockerHub</p> <p>Optimizations upstreamed to main Python trunk</p> <p>Commercial support through Intel® Parallel Studio XE 2017</p>
Intel® Architecture Platforms		
Operating System: Windows*, Linux*, MacOS <sup>1*</sup>		



<sup>1</sup>Intel® Math Kernel Library

<sup>2</sup>Intel® Data Analytics Acceleration Library

## Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.

\*Other names and brands may be claimed as the property of others.

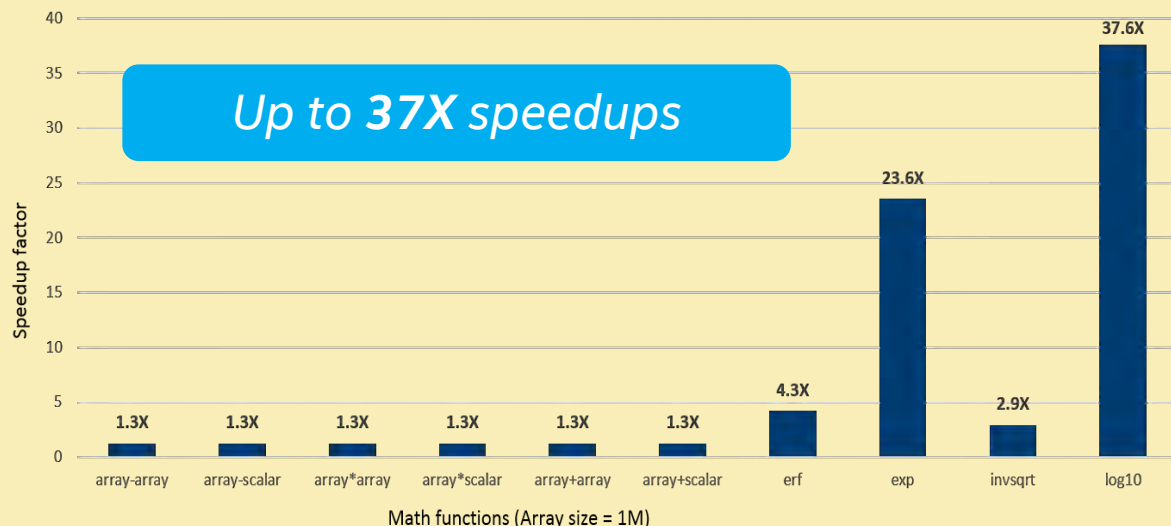
<sup>1</sup> Available only in Intel® Parallel Studio Composer Edition.

# UMath Optimizations & Vectorization to Utilize Multiple Cores, Memory Management

## Intel® Core™ i7 Processor

Intel® Distribution for Python\* Performance Speedups for Select Math Functions on Intel® Core™ i7 Processors

■ Speedup with Intel Python vs pip/numpy



**Hardware:** Intel® Core™ i7-7567U CPU@3.50GHz (1 socket, 2 cores per socket, 2 threads per core), 32GB DDR4 @ 2133MHz. Intel® Xeon® CPU E5-2699 v4@2.20GHz (2 sockets, 22 cores per socket, 1 thread per core-HT is off), 256GB DDR4@2400MHz. Intel® Xeon Phi™ CPU 7250@1.40GHz (1 socket, 68 cores per socket, 4 threads per core), 192GB DDR4 @1200MHz, 16GB MCDRAM@7200MHz in cache mode

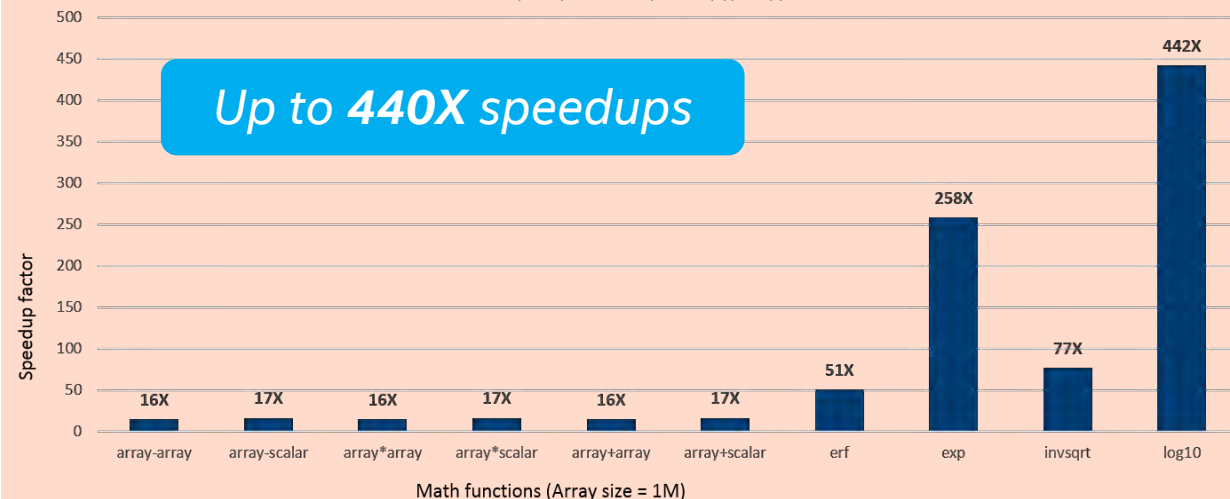
**Software:** Stock: CentOS Linux release 7.3.1611 (Core), python 3.6.2, pip 9.0.1, numpy 1.13.1, scipy 0.19.1, scikit-learn 0.19.0

Intel® Distribution for Python 2018 Gold packages: mkl 2018.0.0 intel\_4, daal 2018.0.0.20170814, numpy 1.13.1 py36\_intel\_15, openmp 2018.0.0 intel\_7, scipy 0.19.1 np113py36\_intel\_11, scikit-learn 0.18.2 np113py36\_intel\_3

## Intel® Xeon™ Processor

Intel® Distribution for Python\* Performance Speedups for Select Math Functions on Intel® Xeon™ Processors

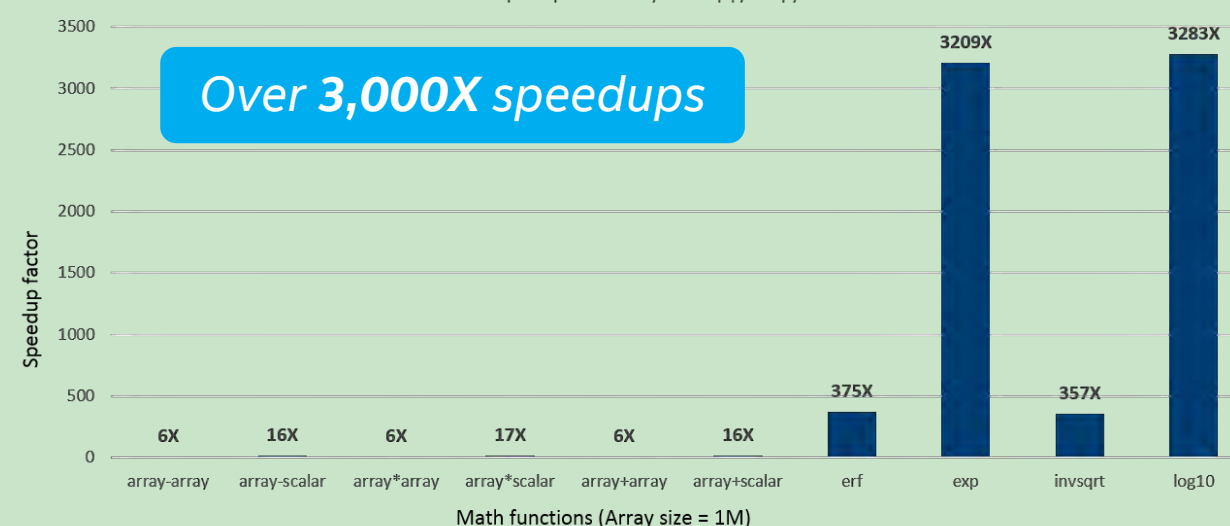
■ Speedup with Intel Python vs pip/numpy



## Intel® Xeon® Phi™ Processor

Intel® Distribution for Python\* Performance Speedups for Select Math Functions on Intel® Xeon Phi™ Processor Family

■ Speedup with Intel Python vs pip/numpy



Software & workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark & MobileMark, are measured using specific computer systems, components, software, operations & functions. Any change to any of those factors may cause the results to vary. You should consult other information & performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more information go to <http://www.intel.com/performance>.

### Optimization Notice

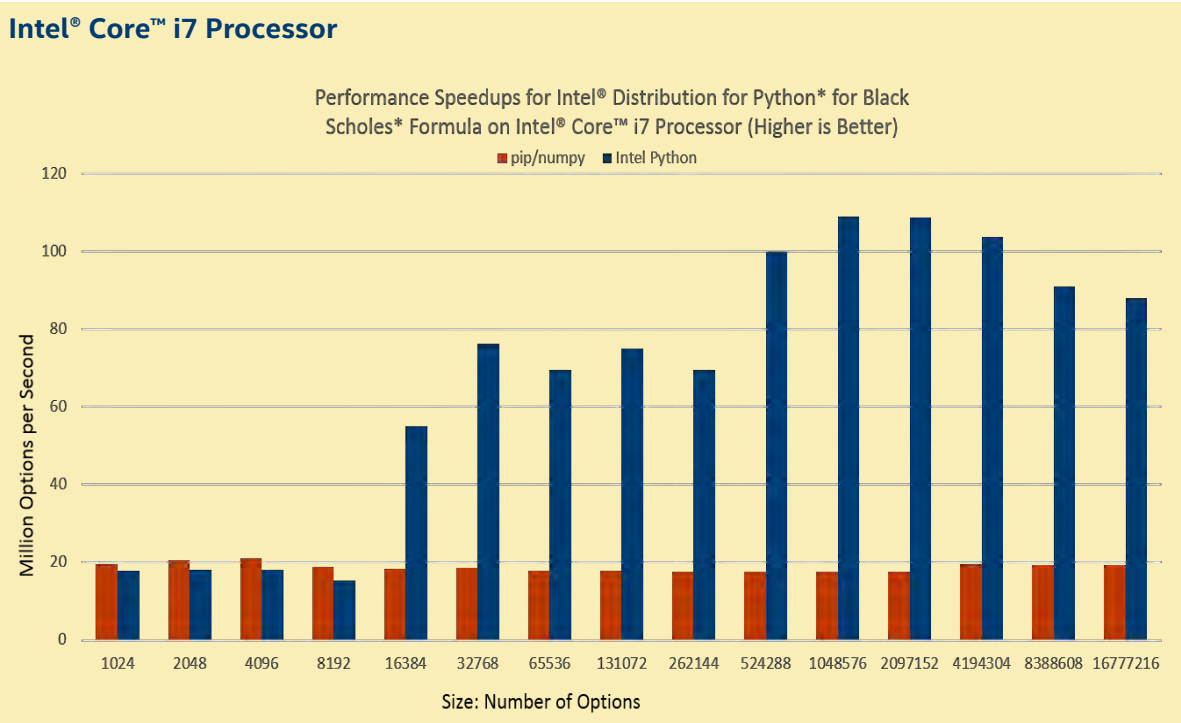
Copyright © 2017, Intel Corporation. All rights reserved.

\*Other names and brands may be claimed as the property of others.

Optimization Notice: Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.



# Performance Speedups for Black Scholes Formula

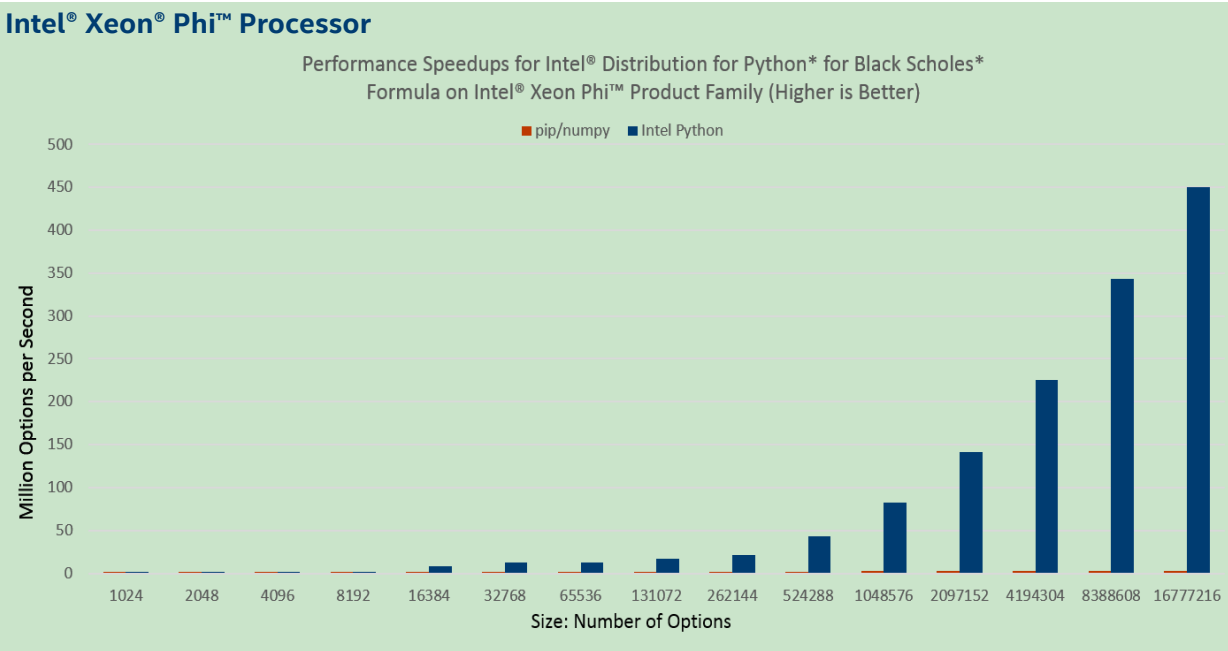
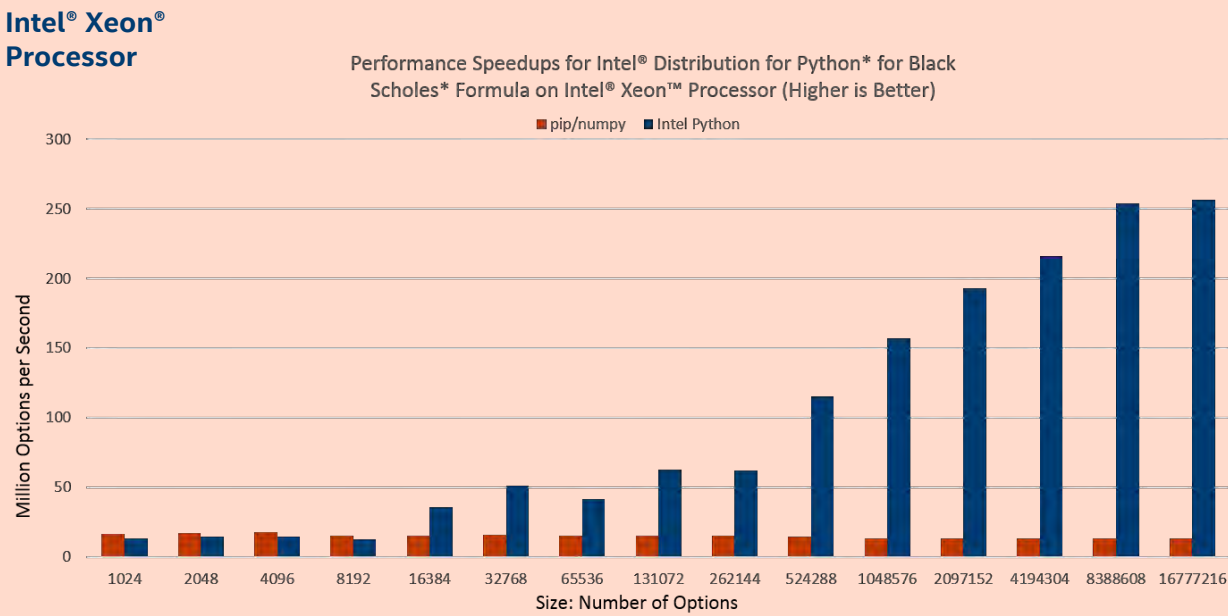


**Hardware:** Intel® Core™ i7-7567U CPU@3.50GHz (1 socket, 2 cores per socket, 2 threads per core), 32GB DDR4 @ 2133MHz. Intel® Xeon® CPU E5-2699 v4@2.20GHz (2 sockets, 22 cores per socket, 1 thread per core-HT is off), 256GB DDR4@2400MHz. Intel® Xeon Phi™ CPU 7250@1.40GHz (1 socket, 68 cores per socket, 4 threads per core), 192GB DDR4 @1200MHz, 16GB MCDRAM@7200MHz in cache mode

**Software:** Stock: CentOS Linux release 7.3.1611 (Core), python 3.6.2, pip 9.0.1, numpy 1.13.1, scipy 0.19.1, scikit-learn 0.19.0

Intel® Distribution for Python 2018 Gold packages: mkl 2018.0.0 intel\_4, daal 2018.0.0.20170814, numpy 1.13.1 py36\_intel\_15, openmp 2018.0.0 intel\_7, scipy 0.19.1 np113py36\_intel\_11, scikit-learn 0.18.2 np113py36\_intel\_3

Software & workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark & MobileMark, are measured using specific computer systems, components, software, operations & functions. Any change to any of those factors may cause the results to vary. You should consult other information & performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more information go to <http://www.intel.com/performance>.



## Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.  
\*Other names and brands may be claimed as the property of others.

Optimization Notice: Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.



# Installing Intel® Distribution for Python\* 2018

## Standalone Installer

Download full installer from  
<https://software.intel.com/en-us/intel-distribution-for-python>

## Anaconda.org

[Anaconda.org/intel channel](https://anaconda.org/intel)

```
> conda config --add channels intel  
> conda install intelpython3_full  
> conda install intelpython3_core
```

## Docker Hub

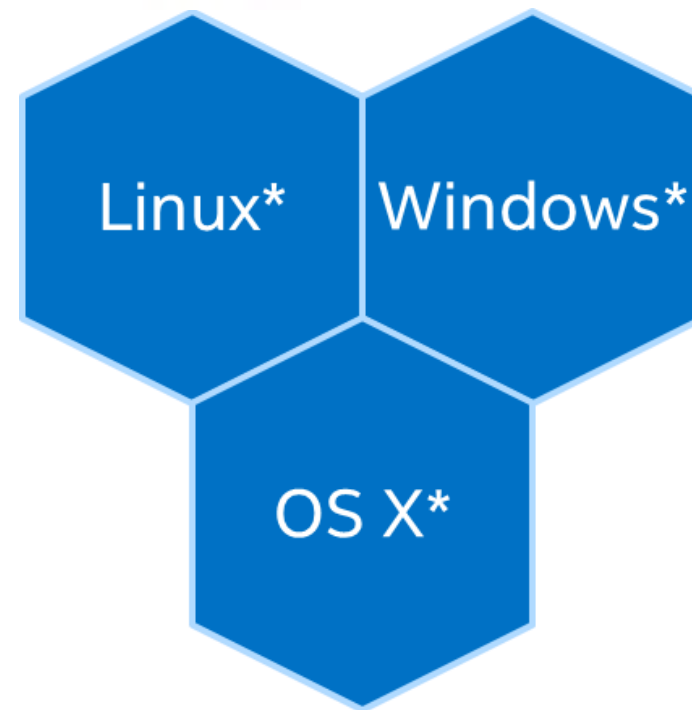
```
docker pull intelpython/intelpython3_full
```

## YUM/APT

Access for yum/apt:  
<https://software.intel.com/en-us/articles/installing-intel-free-libraries-and-python>



2.7 & 3.6





# But Wait.....There's More!



Outside of optimized Python\*, how efficient is your Python/C/C++ application code?



Are there any non-obvious sources of performance loss?



Performance analysis gives the answer!

# Tune Python\* + Native Code for Better Performance

Analyze Performance with Intel® VTune™ Amplifier (available in Intel® Parallel Studio XE)

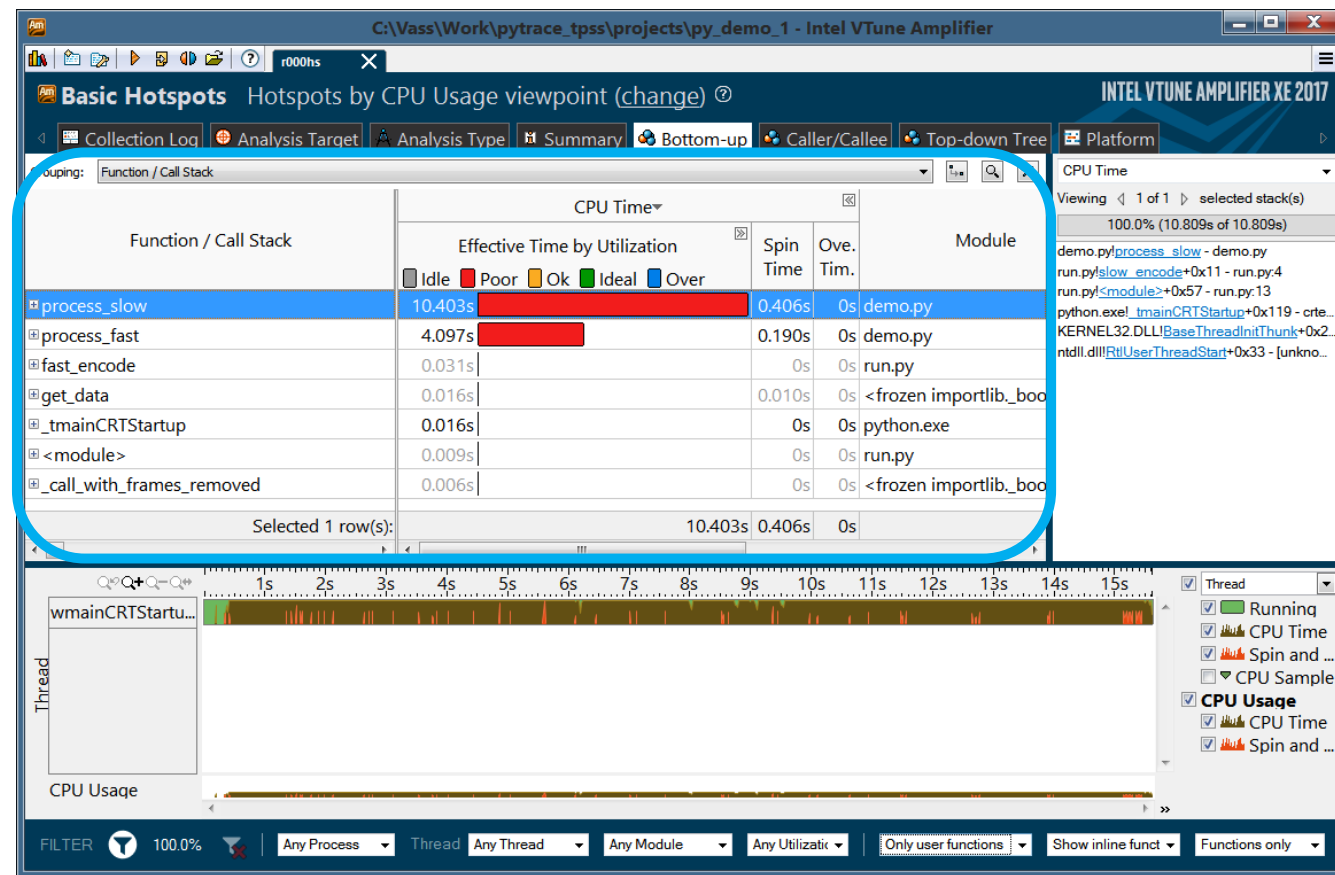
## Challenge

- Single tool that profiles Python + native mixed code applications

- Detection of inefficient runtime execution

## Solution

- Auto-detect mixed Python/C/C++ code & extensions
- Accurately identify performance hotspots at line-level
- Low overhead, attach/detach to running application
- Focus your tuning efforts for most impact on performance



Auto detection & performance analysis of Python & native functions

Available in Intel® VTune™ Amplifier & Intel® Parallel Studio XE

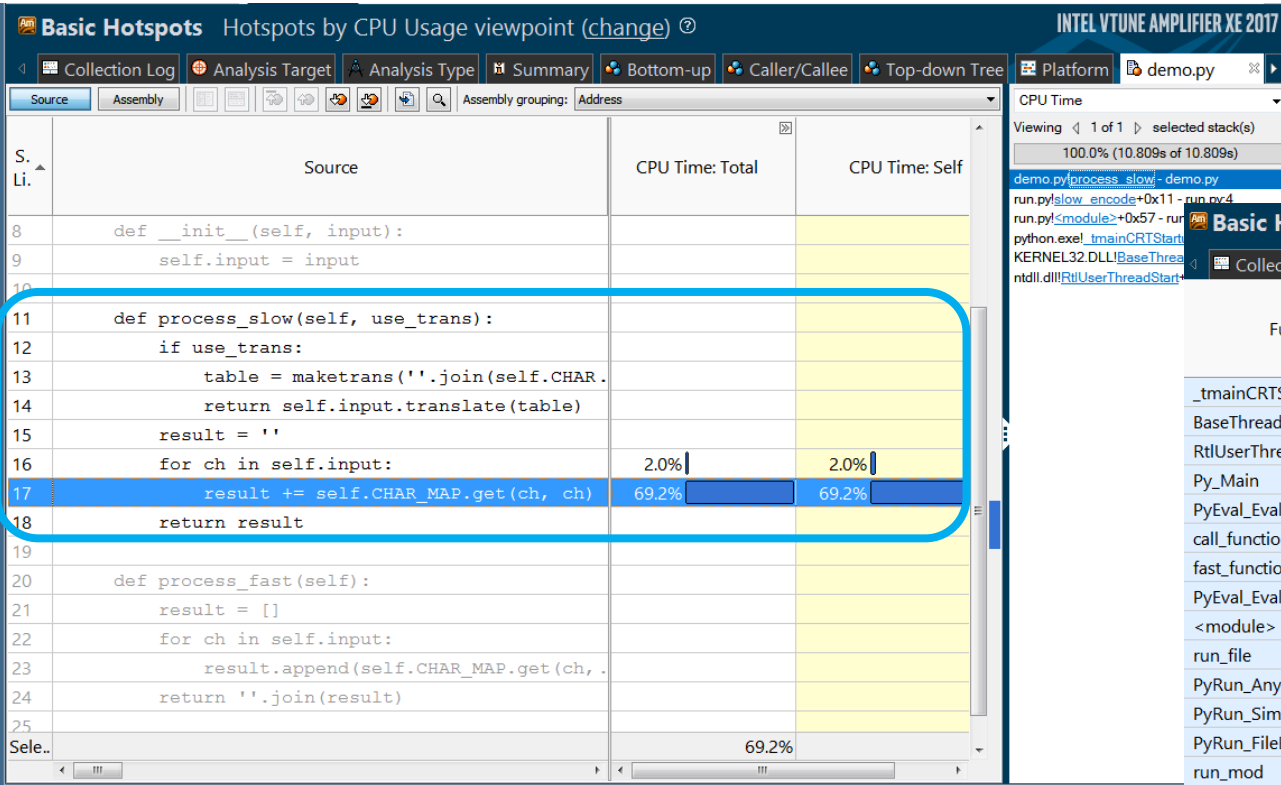
### Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.

\*Other names and brands may be claimed as the property of others.

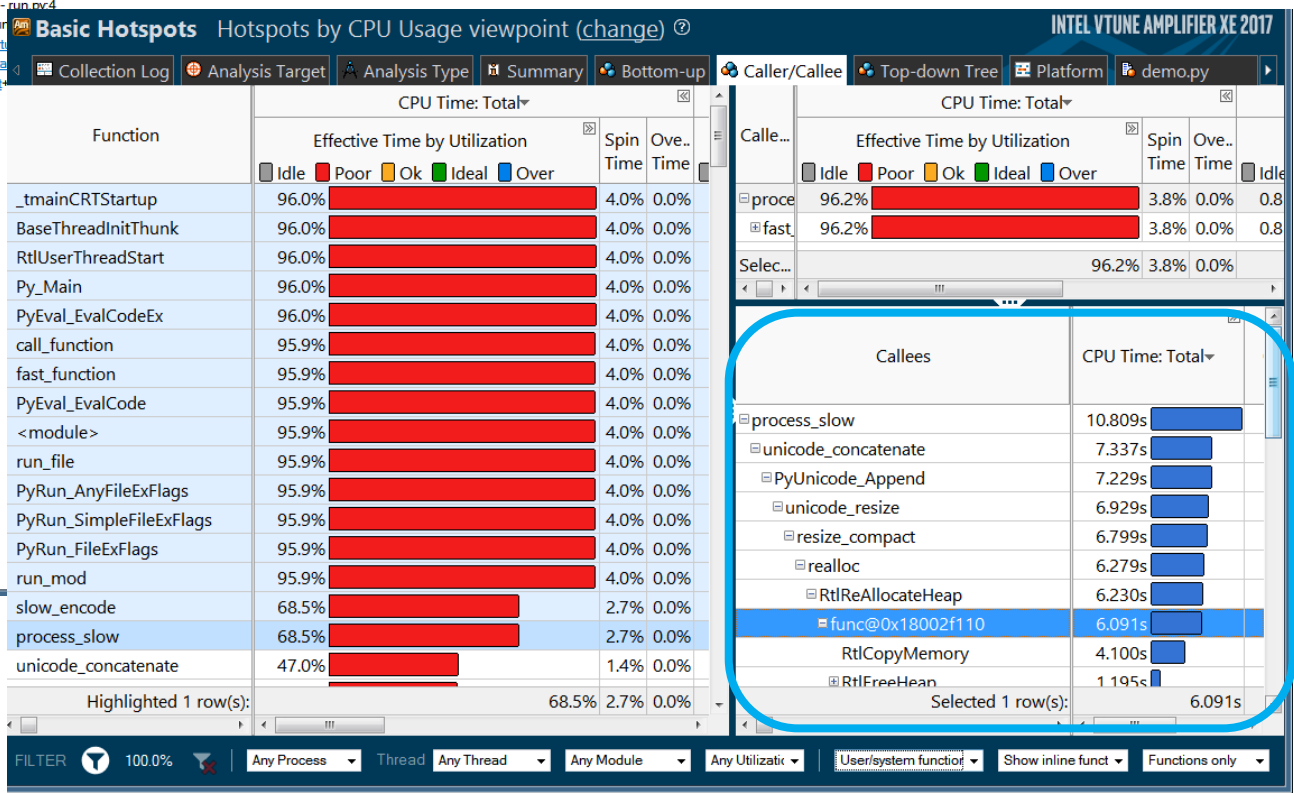


# Diagnose Problem code quickly & accurately



Identifies exact line of code that is a bottleneck

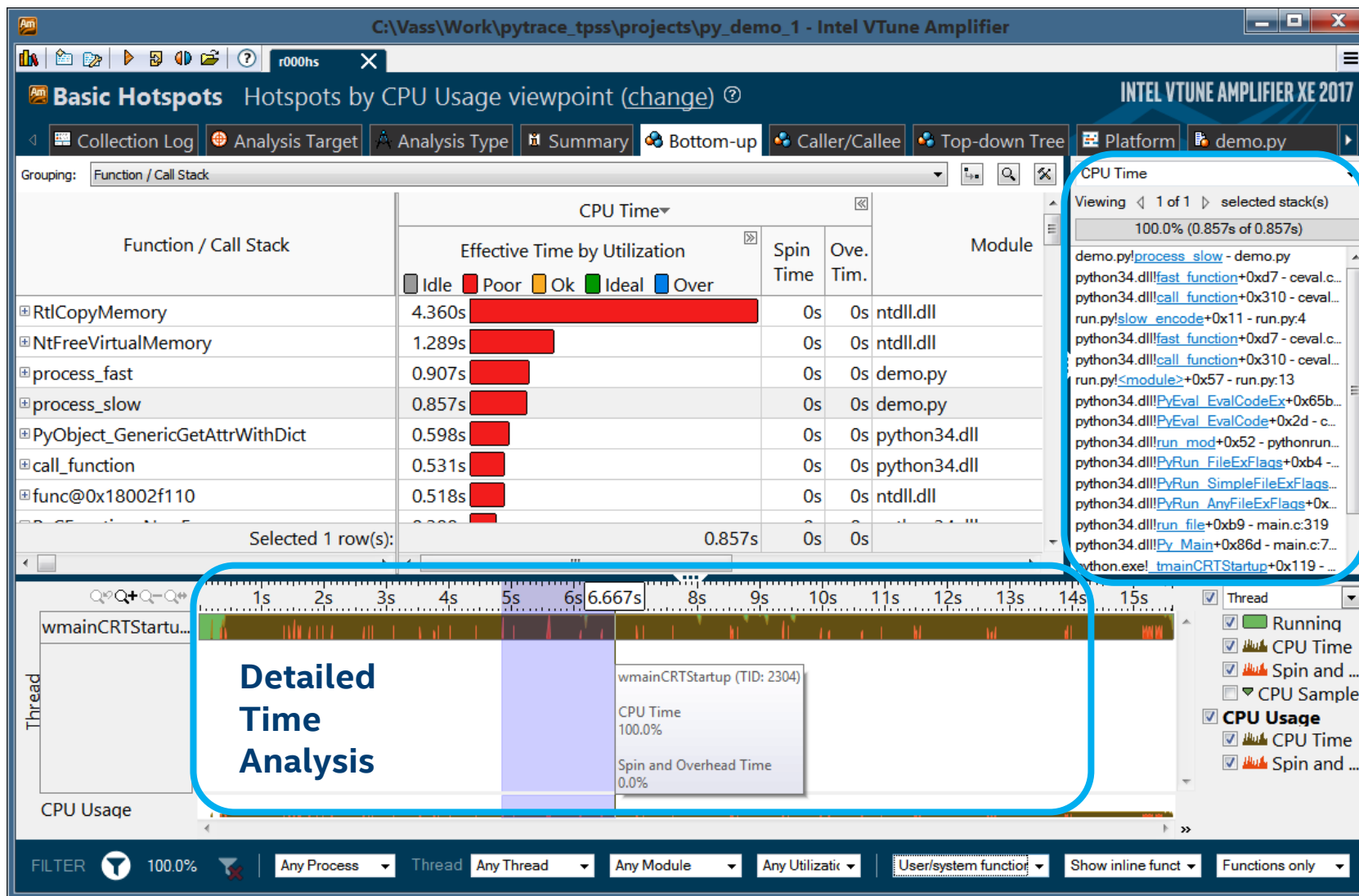
## Details Python\* calling into native functions



### Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.  
\*Other names and brands may be claimed as the property of others.

# Deeper Analysis for Better Insight



Call Stack Listing  
for Python\* &  
Native Code

## Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.  
\*Other names and brands may be claimed as the property of others.





# A 2-prong approach for Faster Python\* Performance

## High Performance Python Distribution + Performance Profiling

### Step 1: Use Intel® Distribution for Python

- Leverage optimized native libraries for performance
- Drop-in replacement for your current Python - no code changes required
- Optimized for multi-core and latest Intel processors

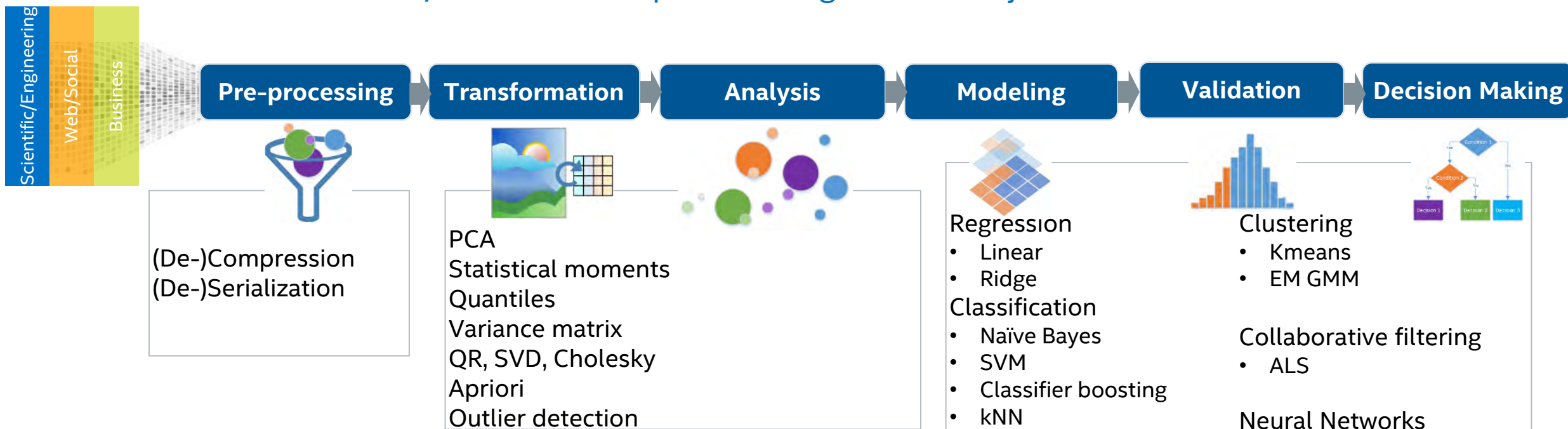
### Step 2: Use Intel® VTune™ Amplifier for profiling

- Get detailed summary of entire application execution profile
- Auto-detects & profiles Python/C/C++ mixed code & extensions with low overhead
- Accurately detect hotspots - line level analysis helps you make smart optimization decisions fast!
- Available in Intel® Parallel Studio XE Professional & Cluster Edition

# PYDAAL

# Intel® Data Analytics Acceleration Library (Intel® DAAL)

- Targets both data centers (Intel® Xeon® and Intel® Xeon Phi™) and edge-devices (Intel® Atom)
- Perform analysis close to data source (sensor/client/server) to optimize response latency, decrease network bandwidth utilization, and maximize security
- Offload data to server/cluster for complex and large-scale analytics



## Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.

\*Other names and brands may be claimed as the property of others.



# Computational Aspects of Big Data

## Volume

- Distributed across different nodes/devices
- Huge data size not fitting into node/device memory

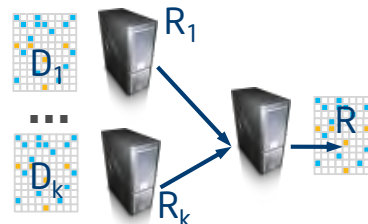
## Variety

- Non-homogeneous data
- Sparse/Missing/Noisy data

## Velocity

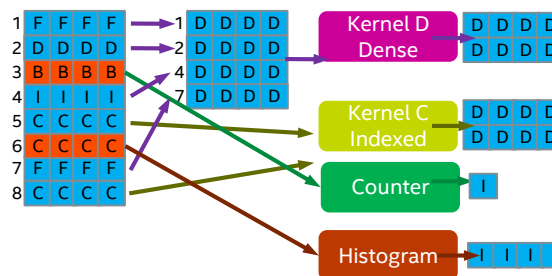
- Data coming in time

### Distributed Computing

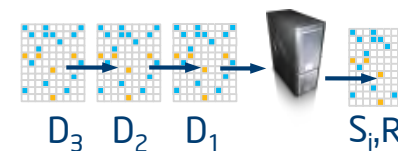


$$R = F(R_1, \dots, R_k)$$

### Converts, Indexing, Repacking

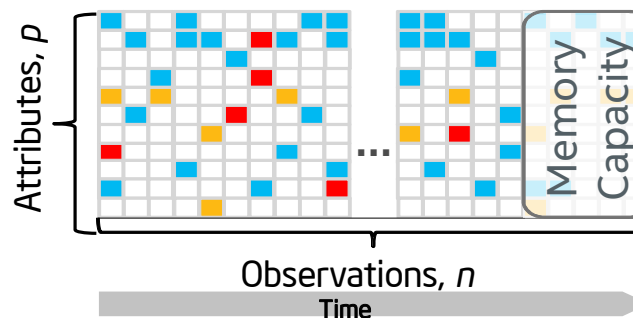
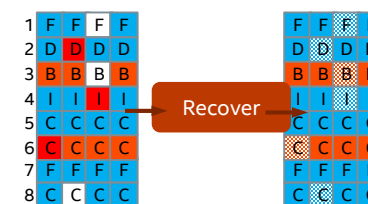


### Online Computing



$$S_{i+1} = T(S_i, D_i); R_{i+1} = F(S_{i+1})$$

### Data Recovery





# Regression

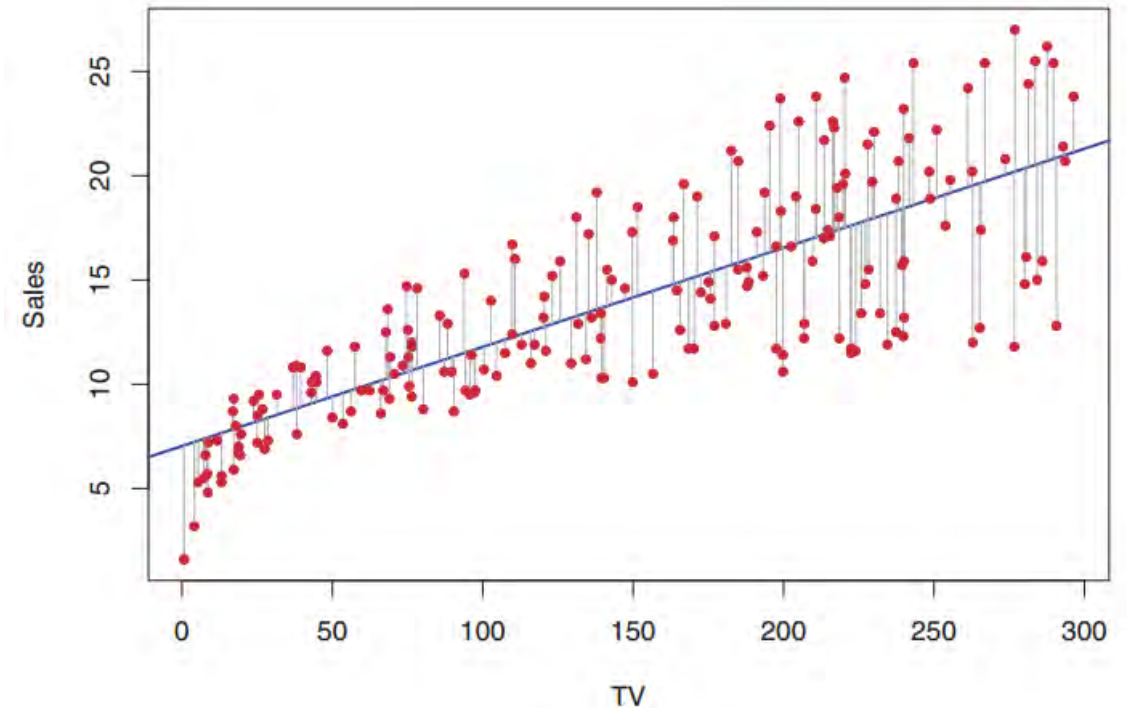
## Problems

- A company wants to define the impact of the pricing changes on the number of product sales
- A biologist wants to define the relationships between body size, shape, anatomy and behavior of the organism

## Solution: Linear Regression

- A linear model for relationship between features and the response

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_N x_N$$



Source: Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani. (2014). *An Introduction to Statistical Learning*. Springer

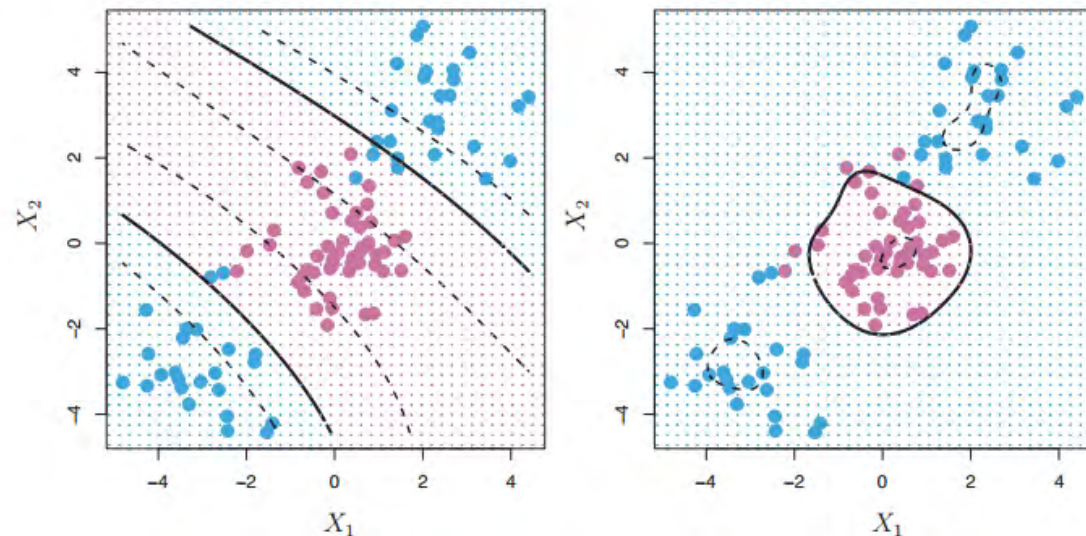
# Classification

## Problems

- An emailing service provider wants to build a spam filter for the customers
- A postal service wants to implement handwritten address interpretation

## Solution: Support Vector Machine (SVM)

- Works well for non-linear decision boundary
- Two kernel functions are provided:
  - Linear kernel
  - Gaussian kernel (RBF)
- Multi-class classifier
  - One-vs-One



Source: Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani. (2014). *An Introduction to Statistical Learning*. Springer

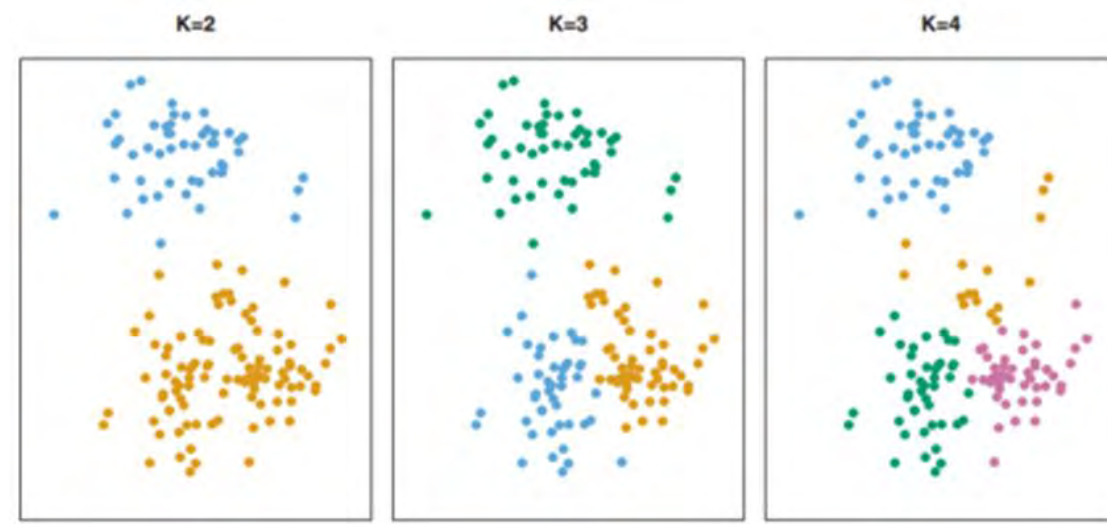
# Cluster Analysis

## Problems

- A news provider wants to group the news with similar headlines in the same section
- Humans with similar genetic pattern are grouped together to identify correlation with a specific disease

## Solution: K-Means

- Pick  $k$  centroids
- Repeat until converge:
  - Assign data points to the closest centroid
  - Re-calculate centroids as the mean of all points in the current cluster
  - Re-assign data points to the closest centroid



Source: Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani. (2014). *An Introduction to Statistical Learning*. Springer

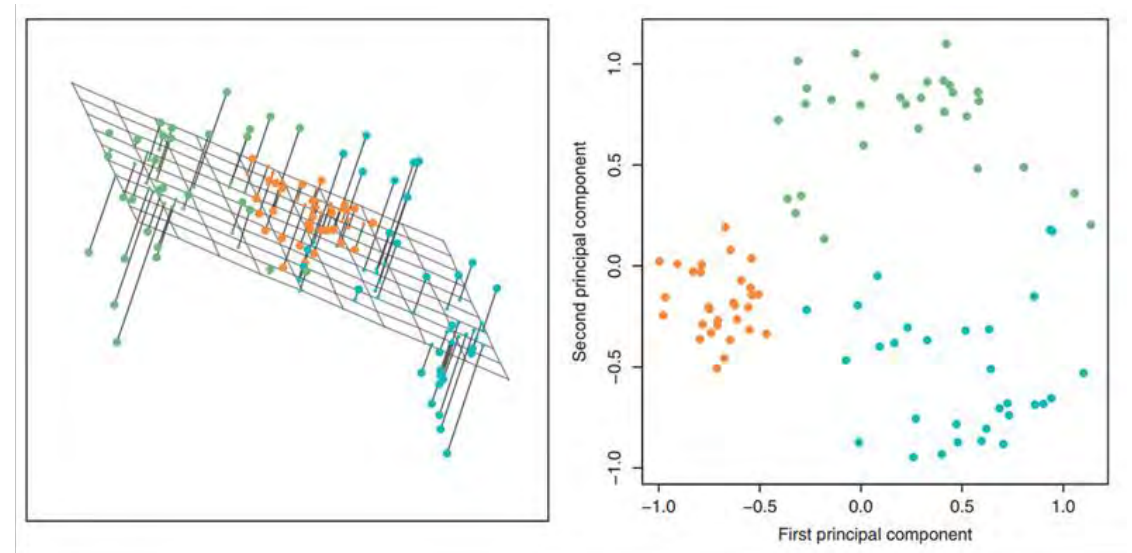
# Dimensionality Reduction

## Problems

- Data scientist wants to visualize a multi-dimensional data set
- A classifier built on the whole data set tends to overfit

## Solution: Principal Component Analysis

- Compute eigen decomposition on the correlation matrix
- Apply the largest eigenvectors to compute the largest principal components that can explain most of variance in original data



Source: Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani. (2014). *An Introduction to Statistical Learning*. Springer



# Demo

[https://github.com/IntelPython/BlackScholes\\_bench](https://github.com/IntelPython/BlackScholes_bench)

<https://github.com/daaltces/pydaal-tutorials>

# CODE THAT PERFORMS AND OUTPERFORMS

Download a *free*, 30-day trial of  
Intel® Parallel Studio XE 2018 today

<https://software.intel.com/en-us/intel-parallel-studio-xe/try-buy>

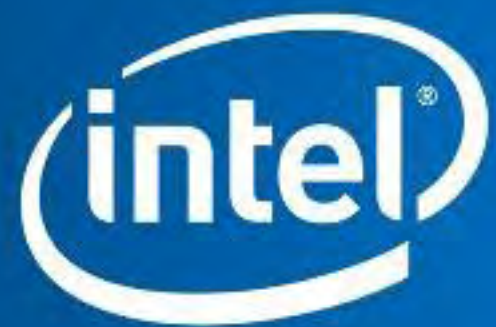
## AND DON'T FORGET...

To check your inbox for the evaluation survey which will be emailed after this presentation.

### P.S.

Everyone who fills out the survey will receive a personalized certificate indicating completion of the training!





Software