Algorithmic Methods of Data Mining AWS Elastic Compute Cloud & Data Science at the Command Line Ioannis Chatzigiannakis Sapienza University of Rome Laboratory 3	 AWS: Elastic Compute Cloud (EC2) AWS EC2 = Elastic Compute Cloud Resizable compute resources in the cloud. Minimizes the time to provision a server. Introduce a new server within minimum delay. Scale capacity up very fast. Quickly modify the capabilities of the compute instance. Introduce additional computational, memory and storage capabilities. Shutdown - or completely remove resources. Scale down very fast. Pay only for the resources you need.
100 \$ 151131 \$ 100	000 (B) (B) (B) (B) (B)
Typical Use Cases Development and Testing Environments Hosting of Databases Data analytics Code repository GPU-assisted machine learning High performance computing Video processing Backup and disaster recovery 	 EC2 Provisioning Options On Demand – Pay for the compute capacity by the hour. No up-front payment or long-term commitment. Short-term, spiky, or unpredictable workloads. Applications development or testing. Spot Instances – Acquire spare capacity up to 90% off the on-demand price. When start/end times are flexible. Applications that are only feasible at very low compute prices. Urgent computing needs for large amounts of additional capacity. Reserved Instances – Significant discount (up to 75%) compared to On-Demand instance pricing. For applications that have steady state or predictable usage. Long term (≥ 1 year) to reduce their total computing costs.

EC2 Instance Types

- General Purpose balance of compute, memory and networking resources.
- Compute Optimized ideal for compute bound applications that benefit from high performance processors.
- Memory Optimized deliver fast performance for workloads that process large data sets in memory.
- Accelerated computing use hardware accelerators, or co-processors, to perform functions, such as floating point number calculations, graphics processing, or data pattern matching, more efficiently than is possible in software running on generic CPUs.
- Storage optimized for workloads that require high, sequential read and write access to very large data sets on local storage.

EC2 Instance Types & Resources

- CPU 64-bit Arm, AMD EPYC 7000, Intel Xeon Platinum 8175M, Intel Xeon E5-2676.
 - 1 192 virtual CPUs 1 thread = 1 vCPU.
- ▶ Memory 1 ... 512 GB.
- Network up to 100 Gbps.
- Storage
 - Amazon Elastic Block Store (EBS) easy to use, high performance block storage service.
 - 0...60 TB NVMe SSD ensure best IOPS (Input/Output operations per second).
- Hardware Accelerators
 - NVIDIA Tesla V100 GPUs, NVIDIA K80 GPUs, NVIDIA T4 Tensor Core GPUs.
 - AWS Inferentia Chips.
 - Xilinx Virtex UltraScale+ VU9P FPGAs



Available OS & Software

Operating Systems

- Linux/Unix Amazon Linux, Debian, Ubuntu, Red Hat, CentOS, SUSE, FreeBSD, Gentoo, Mint, ...
- Windows Server 2019, Server 2016, Server 2012.
- Databases PostgreSQL, MySQL, MongoDB, Neo4J, Oracle Enterprise, Microsoft SQL, ...
- AWS Marketplace a wide selection of commercial and free software from well-known vendors.

Pricing Examples

- General Purpose
 - t2.micro Linux or Windows 2 vCPUs + 4 GB 750 hours free per month, \$0.05/h
 - a1.xlarge Linux 4 64-bit ARM vCPUs + 8 GB \$0.1152/h
 - a1.xlarge Linux 4 64-bit ARM vCPUs + 8 GB \$0.1152/h
 - m5.24xlarge Linux 96 Xeon vCPUs + 337 GB \$5.136/h
 - m5.24xlarge Windows 96 Xeon vCPUs + 337 GB \$9.552/h
- Compute Optimized
 - c5.xlarge Linux 4 Xeon vCPUs + 8 GB \$0.192/Hour
 - c5.24xlarge Linux 96 Xeon vCPUs + 192 GB \$4.608/Hour
- Hardware Accelerators
 - p3.2xlarge Linux 1 NVIDIA Tesla V100 GPUs + 8 Xeon vCPUs + 61 GB – \$3.305 per Hour
 - p3dn.24xlarge Linux 8 NVIDIA Tesla V100 GPUs + 96 Xeon vCPUs + 768 GB – \$33.711 per Hour



101 (B) (2) (2) (2) 2 000

10, 00, 12, 12, 12, 2, 000

Amazon Elastic Block Store (EBS)

- Easy to use, high performance block storage service.
- Targeting both throughput and transaction intensive workloads.
 - Can be used for relational and non-relational databases.
 - Enterprise applications.
 - Big data analytics engines.
 - General purpose file systems.
 - Media workflows.
- Highly availability and durability 99.999%
- Virtually unlimited scale as little as a single GB of storage, or scale up to petabytes of data.
- Secure encryption of data at-rest, data in-transit, and all volume backups.

EBS Volume Types - HDD based

- Throughput Optimized HDD (ST1) ideal for frequently
 - accessed, throughput-intensive workloads.
 - Large datasets and large I/O sizes, such as MapReduce, Kafka, log processing, data warehouse, and ETL workloads.
 - Low cost HDD volume.
 - Volume Size: 500 GB 16 TB.
 - Max IOPS/Volume: 500
 - Max Throughput/Volume: 500 MB/s
 - Price: \$0.045/GB-month
- Low-cost HDD (SC1) ideal for less frequently accessed workloads with large, cold datasets.
 - Colder data requiring fewer scans per day.
 - Volume Size: 500 GB 16 TB.
 - Max IOPS/Volume: 250
 - Max Throughput/Volume: 250 MB/s
 - Price: \$0.025/GB-month

EBS Volume Types - SSD based

- Provisioned IOPS SSD (IO1) high performance SSD volume designed for latency-sensitive transactional workloads.
 - I/O-intensive NoSQL & relational databases.
 - Volume Size: 4 GB 16 TB.
 - Max IOPS/Volume: 64,000
 - Max Throughput/Volume: 1,000 MB/s
 - Price: \$0.125/GB-month + \$0.065/provisioned IOPS
- Default EBS volume type (GP2) ideal for suitable for a broad range of transactional workloads.
 - Boot volumes, low-latency interactive apps, dev & test.
 - Volume Size: 1 TB 16 TB.
 - Max IOPS/Volume: 16,000
 - Max Throughput/Volume: 250 MB/s
 - Price: \$0.10/GB-month

Choose Region

A Nochronit Construction And Annual Construction Annual Constructio Annual Construc	ndephotechngket-ex-eard 1	- 9.6 0 0.5m		4 N D 6 K * 5
v				
AWS Manageme	nt Console			
#M5 services			Stey connected to the ga	
Field Services Transmission server, Reparate a serveryon Q, Exemption Relationed Databases Service data	Indone, RCK		e man	ARTIST PARADO SHOT
 Insertly risked services R ICI 	R Date Colore Server	0.000	Explore ANS	All Participated to restrict C All Participation (control of f All Participation) to restrict (
2.0	Anana Inninan	D und lange lage the	Feer Digital Training	
O erten	El Anazorito	25 men	fort assessments 200+ and HMS products and new	
 Art Operation Art 	 HX Initial descentions 	 By send the second secon	Amazan Lapitaine Set hank ar atticke	
All service Strength Kz spread Q service	 Jacobia onuncement vir disament Technologies 	 Security, Identity, & Complexes MM Research Access Namegar Comm 	WYS-Cartification Depict the summary processing contribution	Europe (Hand Account) Europe (Hand Account) Europe (Hand Account) Europe (Hand Account)
kent Earth Brontolt	Anarolisist	Gently Nanapr Gamblers	Aneoer 015 for AVS	
Lar antinu. Sppiluation Reportany MRT-Colourity	(C) Hanagement & Decembers AND Constructions	Imperiate Resident Marine	benotify the iteration for arrite and paraterisings	South among Saturdade accord.

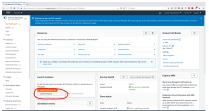


(D) (Ø) (2) (2) (2) 2 (0)

Open EC2 Service Dashboard



Launch Instance



https://www.f.samala.acc.anasocom/wij/U/Jona/reporter and 18.autilitation



Select Amazon Machine Image (AMI)

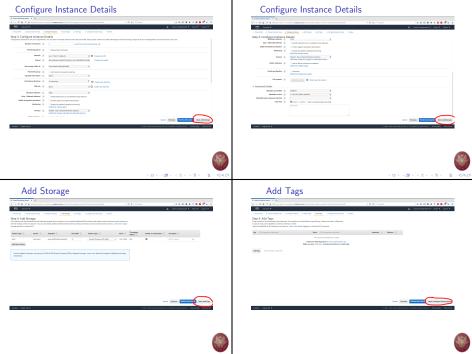
- C & 0 A10	Street contraction of the street the street the street the street of the street street the street street the street stree	100.05
En lersten V	4 terres Designments #	Interf. W. Dagert #
mount 100000		
ep 1: Choose an concurs Atta	Amazon Machine Image (AMI) National Control of Control	Consel and Date
Testining ()	Annual Technologies (1997) Annual Technologies Annual Technol	Select SALUE (1995)
	Ind the films place lines i \$1955, 600 Yebsen Type: uni BHC (Valid) (MAR) (644 HBC (2016) (2640	Defend
	REELEN Dispers have 13 99 \$44,000 years by a since State Colling of the State S	Calues Calaba (chil) Calaba (con)
	Reamb Genery (2014) (2014	Const 1000 Const 1000 Const 1000
	A try to longing of data motions ¹ by datase RE. A set to longing of datase motions ¹ by datase RE. In the set of the s	Ride andly deploy meanwind
	Beels Server 16 (H121 (H101, 103 Telever See - on 0000/2001/h131 (H42 on) / on 0000/0010/01 (H42 Am)	5691

Choose Instance Type

100	latin V						A hereit Ontoinenti • hele	dia kaominina
	a Long ages for							
nanel(mixed emotionality year applica		or lypes and how they care	er chui annes fui sa se eri pur angulig ceis	oppination. They have saying consid	nations of DN, reamony strengt, and set	undairg capacity, and give you like first	hilly to absorb th
Connector	undentanik Corriente (Harrishine KOA) Familiy	1 v2NA 23.0rs, Intel Bound	weight bill marriery, Ditt a	Barray (185)	mitters through \$45 (2)	In particulation ()	Induct Parlomance (2)	Pathque
	levent payme	27479		- 65	011110		Los to Roberton	140
•	lenest payme	ti misu			05 140		Line to Roberton	50
	limest payme	15 small			01110		Los to Nobeleville	54
	lanari payna	threadware.		4	085.000		Line to BODOWN	740
	licitani parprise	2018			001140		L2+ 3/ 800/03/	740
	factorial purpose	0.00		10	083 199		Million	740
	tenent papers	0.5kmp			000 149		Maderare	744
	tenent autore	Daneo		0	000 149	14	1010 9 Diplot	144
	tenent autore	Claiming .			001 +45	144	Up to 3 Digital	144
	Denieral guargione	tia seat			101.00	-	Up to 2 Equity	ten.
	Denient gurgener	the resolute		4	101.00	54	Same Paper	_
						Canvel Previous	derive and Leased Name Lands	sen inniaense žariai



000 5 (5) (5) (5) (0)



9.00 S (S)(S)(S)(S) (S)

Configure Security Group

o d' & 0 & Hochevert	onen an			4 N D 6 K 7 2 1	a Pierre
😂 levies 🖲				A land Designments * Intered *	have a
Chose and Echoaccases for	Lindgenomic Lindbing Li	Artigo Anterportunity Street			
unresident access to the HTP and HTPD		page per consolisi sino institue specific inefficia en ar admititum an existing are better, Leare mare ab	nie para innianze Are mannyle, if pro most in set op a anto od timolon 112 anuarly yrodyk.	never and allow interest india in wash your induces, addi	ndre find alsor
	Constant elimination prop-				
twowits group-mate					
Description	www.wardcound.tempt	117 34 34 MI-ROOM			
ter G	Pretrained (2)	Pet Baye (3)	laws ()	Description (3)	
601 +	304	10	Gastern + \$305.0		0
Additute					
hannan biyan jin b			\$ SEC. ANA JUNC		n ad anna
					2 040
					2 940

Launch Instance



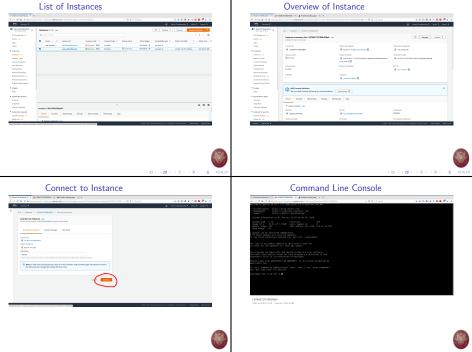


Instance Ready to Launch

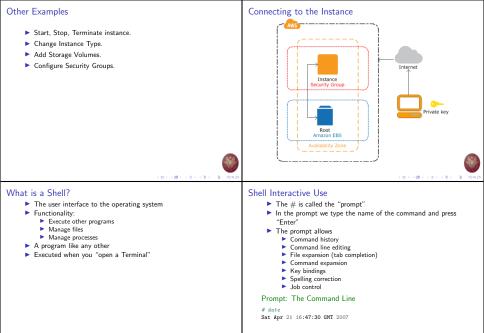








900 S (S) (S) (S) (B)



Error Handling

If we type a wrong command, an error message appears

Prompt: The Command Line

datee

datee: no such file or directory

- The error message states that either the file or the folder (directory) was not found
 - In the prompt all commands are assumed to be connected to a file ...
- ▶ The arrow keys ↑↓ allow to look-up previous commands
- \blacktriangleright The arrow keys $\leftarrow \rightarrow$ allow to move within the same command line

Terminating Command Execution

- We can interrupt the execution of a command by pressing ctrl-c
- We can "freeze" the output of the execution of a command by pressing ctrl-s
 - To "un-freeze" the output of a command we use ctrl-q
 - Note only the output is frozen not the actual execution
- ▶ To close a terminal we use ctrl-d
 - We may need to press multiple times ctrl-q
 - All programs currently running will terminate



Manual Pages

- The command man allows to access the manual pages
- Manual pages are organized in categories
 - 1. Commands Is, cp, grep
 - 2. System Calls fork, exit
 - 3. Libraries
 - 4. I/O Files
 - 5. File Encoding Types
 - 6. Games
 - 7. Miscellaneous
 - 8. Administrator's Commands
 - 9. Documents
- We can request a page from a specific category man [category] [topic]

Manual Pages

FORR(2)	Minix Programmer's I	Manua I	FORK(2)
NAME fork – cre	ate a new process		
SYNOPSIS #include < #include <	(sys∕types.h) unistd.h)		
pid_t fork	(void)		
DESCRIPTION Fork cause is an exac	s creation of a new process. ' t copy of the calling process :	The new process (c except for the foll	hild process) owing:
	hild process has a unique proc	ess ID.	
	hild process has a different ss ID of the parent process).		D (i.e., the
The c	hild process has its own copy	of the parent's	descriptors.

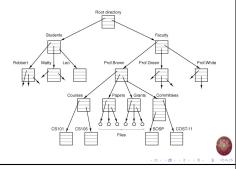


101 (B) (2) (2) (2) 2 000

File System

- All system entities are abstracted as files
 - Folders and files
 - Commands and applications
 - I/O devices
 - Memory
 - Process communication
- The file system is hierarchical
 - Folders and files construct a tree structure
 - The root of the tree is represented using the /
- The actual structure of the tree depends on the distribution of Linux
 - Certain folders and files are standard across all Linux distributions

File System Example



Standard Folders

- /bin Basic commands
- /etc System settings
- /usr Applications and Libraries
- /usr/bin Application commands
- /usr/local Applications installed by the local users
- /sbin Administrator commands
- /var Various system files
- /tmp Temporary files
- /dev Devices
- /boot Files needed to start the system
- /root Administrator's folder

Example of File Metadata

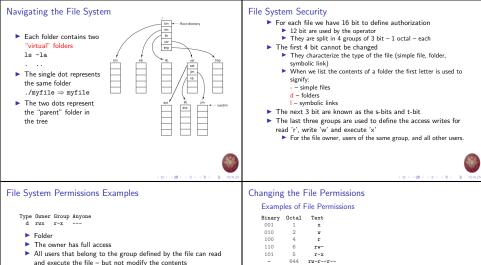
ls -la

lrwxrwxrwx	1	bin	operator	2880	Jun	1	1993	bin
-rr	1	root	operator	448	Jun	1	1993	boot
drwxr-sr-x	2	root	operator	11264	May	11	17:00	dev
drwxr-sr-x	10	root	operator	2560	Jul	8	02:06	etc
drwxrwxrwx	1	bin	bin	7	Jun	1	1993	home
lrwxrwxrwx	1	root	operator	7	Jun	1	1993	lib
drwxr-sr-x	2	root	operator				1992	
drwx	2	root	operator	512	Sep	26	1993	root
drwxr-sr-x	2	bin	operator		Jun			
drwxrwxrwx	6	root	operator	732	Jul	8	19:23	tmp
drwxr-xr-x	27	bin	bin				1993	
drwxr-sr-x	10	root	operator	512	Jul	23	1992	var



100 5 15 15 15 10 000



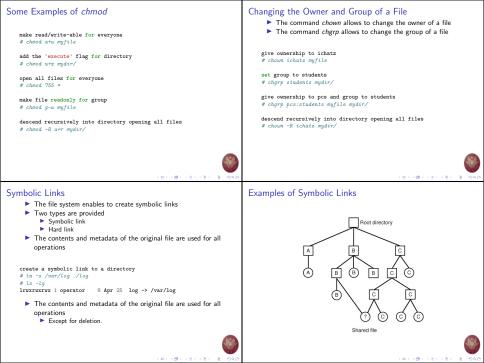


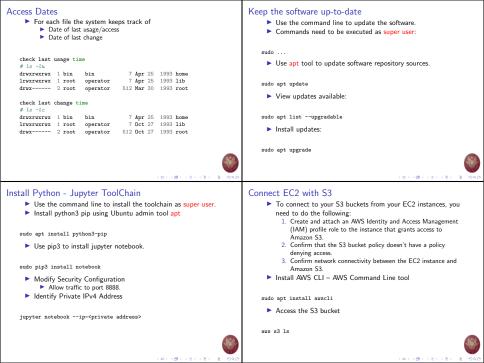
- All other users cannot access the file or execute it
- To access a folder we use the command *cd* given that we have
- Profaccess a router we use the command cd given that we have permission to execute 'x'

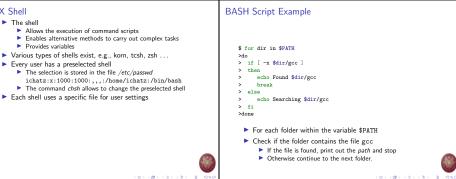
- The command chmod allows to modify the permissions
- There are 2 way to define the new permissions
 - 1. Defining the 3 Octal e.g., 644
 - 2. By using text e.g., a+r



10, 10, 12, 12, 12, 2, 000







Command line

UNIX Shell

hash

bash-4.4.20#

- Left part of # can be changed.
- Right part of # is used to type in commands.
- Offers certain built-in commands
 - Implemented within the BASH source code
 - These commands are executed within the BASH process
- Allows to execute scripts
 - For this reason it is called a UNIX programming environment

Built-in Commands

	nanus	
Command	Description	Exception
cd	Change Folder	cd
declare	Set a variable	declare myvar
echo	Print out a text to the standard out-	echo hello
	put	
exec	Replace bash with another process	exec 1s
exit	Terminate shell process	exit
export	Set a global variable	export myvar=1
history	List of command history	history
kill	Send a message to a process	kill 1121
let	Evaluate an arithmetic expression	let myvar=3+5



Built-in Commands

Command	Description	Exception
local	Declare a local variable	local myvar=5
pwd	The current folder	pwd
read	Read a value from standard input	read myvar
readonly	Lock the contents of a variable	readonly myvar
return	Complete a function call and return a	return 1
	value	
set	List declared variables	set
shift	Shifts the command parameters	shift 2
test	Evaluate an expression	test -d temp
trap	Monitor a signal	trap "echo Signal" 3

UNIX Pipes

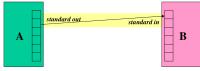
 General idea: The input of one program is the output of the other, and vice versa.



Both programs run at the same time.

UNIX Pipes

Often, only one end of the pipe is used.



This can be done using intermediate files.

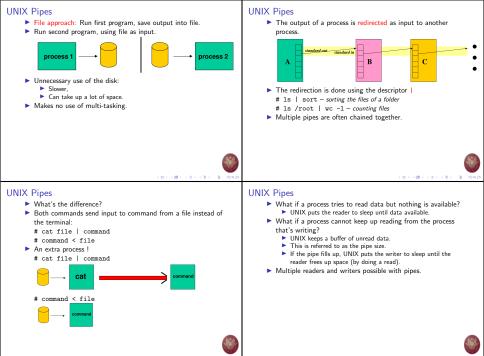
UNIX Pipes

- Commands produce an output using the descriptor > the output is redirected to a file
 - # ls > filelist
- A new file is created under the name filelist
- If the file already exists, the new file will replace the old one.
- We can use the descriptor >> to redirect the output to an existing file
 - # ls -lt /root/doc >> /root/filelist
- The commands that require input using the descriptor < the input is redirected from a file</p>
 - # sort < /root/filelist</pre>

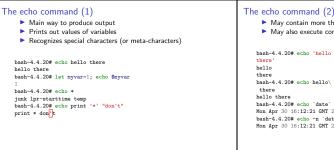


1011001001001000000

100 5 15 15 15 10 000



UNIX Pipes • Examples of filters: • Sort • Input: lines from a file. • Output: lines from a file. • Output: lines from a file. • Output: lines that match the argument. • Sed • Programmable stream editor.	 Processes We may execute commands in series by using the delimeter ; Commands are executed one by one. When the first is completed, the next one starts. When the last command is completed, we get a new prompt # who sort ; date We may execute commands in the background using the delimeter & The commands are executed and a new prompt is provided immediately # pr junk 1pr & The execution of a command results to a new process The command ps shows up in the list of active processes The command wait is active until all the commands executed using the delimeter k complete.
List of processes	Process management
<pre>// ps -a PID TTY THE CMD 106 c1 0:01 -ah 4114 c0 0:00 -bh 214 c0 0:00 -ah 6762 c1 0:00 pa -a 6762 c1 0:00 getty 90 c3 0:00 getty 90 c3 0:00 getty • Parameter a - list all the commands created by consoles • Column TPD - unique ID of the process • Column TIME - total execution time • Column CMD - the name of the command</pre>	 To terminate aprocess we use the command kill [PID] We may change the priority of a process prefix nice # nice pr junk lpr & We may delay the execution of a command prefix at # at 1500 ls -1 / /rooot /dir wc > allfiles pr allfiles lpr ; date > lpr-endtime & date > lpr=starttime To at: /usr/spool/at/07.111.1500.67 created #
l	



The echo command (2)

- May contain more than 1 lines
- May also execute commands

Mon Apr 30 16:12:21 GMT 2007 bash-4.4.20# echo -n `date` " " Mon Apr 30 16:12:21 GMT 2007 bash-4.4.20#



Meta-characters

- The character ? defines any single character, e.g., ls /etc/rc.????
- The character * defines multiple characters, e.g., ls /etc/rc.*
- The array [...] defines a specific set of characters, e.g. ls [abc].c
- The use of the above meta-characters is also called filename. substitution
- We may use these meta-characters in any combination within command execution
- The following command is disabled
 - mv *.x *.y

Shell Variables

- The shell allows the declaration of variables.
- Initial values of variables are defined in the user settings file
- The scope of the variables is connected with the session Or until the user removes them
- The variables with UPPER-case letters are global they are transfered to all processes executed by the shell
- The variables with LOWER-case letters are local they are accessible only by the shell process

		ME
t	e	rm

The path to your home directory # The terminal type



101 (#112) (2) (2) 2 OQ

Shell Variables

- We may use variables at the command line
- We use the descriptor \$

```
bash-4.4.20# myvar="hello"; echo $myvar
hello
bash-4.4.20# myvar="ls -la"
bash-4.4.20# $myvar
lrvxrvxrvx 1 bin operator 2880 Jun 1 1993 bin
-r-r-r-r 1 root operator 448 Jun 1 1993 boot
drvxr-sr-x 2 root operator 11264 May 11 17:00 dev
```

Special Variables

Some special variables are provided

Variable	Description
USER	User name
HOME	Home folder of user
TERM	Type of terminal
SHELL	Name of shell
PATH	List of folders to look for commands
MANPATH	List of folders to look for manual
	pages
PWD	Active folder
OLDPWD	Previously active folder
HOSTNAME	Name of the system



Variable Handling

- The commands env, printenv provide a list of GLOBAL variables
- The command set provides a list of LOCAL variables
- To declare a new GLOBAL variable we use the command export
- Variable type is define by content type
 - String variables myvar = "value"
 - Integer variables declare -i myvar
 - Constant variables readonly me="ichatz"
 - Array variables declare -a MYARRAY MYARRAY[0]="one"; MYARRAY[1]=5; echo \${MYARRAY[*]}
- The names of the variables are case-sensitive
- The command unset removes a variable

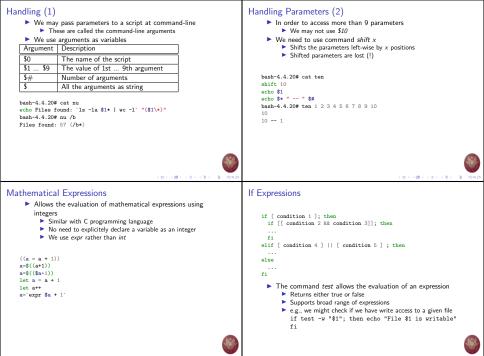
Creation of scripts

- Scripts are used as if they were commands/applications
 Defined by a source file
- We execute the script using the command sh
 - Or directly by setting execute access permissions

```
bash-4.4.20# echo 'who | wc -l' > nu
bash-4.4.20# cat nu
who | wc -l
bash-4.4.20# sh nu
1
bash-4.4.20# chmod a+x nu
bash-4.4.20# nu
1
```



・ロト (使) (注) (さ) (さ) (さ)



Evaluation using test

Expression	Description
-gt	Greater or equal
-ge	Greater
-lt	Smaller
-le	Smaller or equal
-eg	Equal
-ne	Not Equal
-n str	Size of the string bigger than 0
-z str	Empty string
-d file	The file is a folder
-s file	A non empty file
-f file	The file exists
-r file	Read access to file
-w file	Write access to file
-x file	Execution access to file

Evaluation Example (1)

```
elif [ ! -r "$filename" ]; then
    echo "File is not readable"
    exit 1
fi
```

. .

101101101121121 2 000



Evaluation Example (2)

bash-4.4.20# cat check.sh
#!/bin/bash
TMPFILE = "diff.out"

diff \$1 \$2 > \$TMPFILE

if [! -s "\$TMPFILE"]; then
 echo "Files are the same"

else

more \$TMPFILE

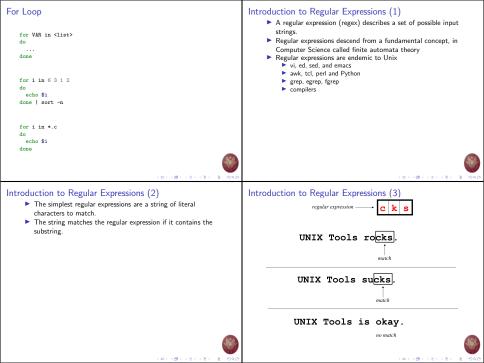
fi

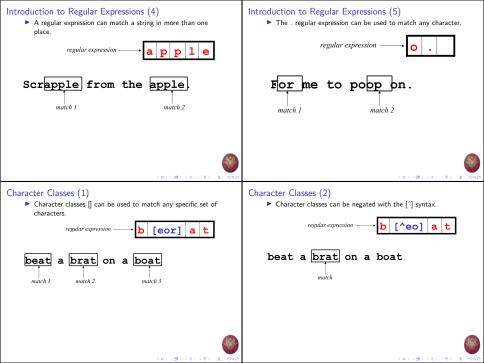
```
if [ -f "$TMPFILE" ]; then
  rm -rf $TMPFILE
fi
```

Boolean expressions

```
if [ condition 1 && condition a]; then
    if [ condition 2 || condition b]; then
    ...
    fi
    elif [ ! condition 3 ] ; then
    ...
    fi
```







Character Classes (3)

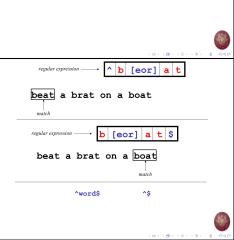
- [aeiou] will match any of the characters a, e, i, o, or u
- [kK]orn will match korn or Korn
- Ranges can also be specified in character classes
- ▶ [1 9] is the same as [123456789]
- ▶ [abcde] is equivalent to [a − e]
- You can also combine multiple ranges
- ▶ [abcde123456789] is equivalent to [a e1 9]
- Note that the character has a special meaning in a character class but only if it is used within a range,
- ▶ [-123] would match the characters -, 1, 2, or 3

Named Character Classes

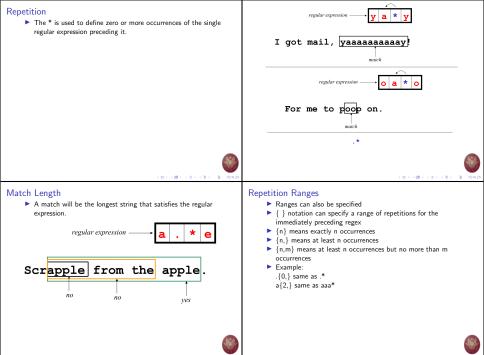
- Commonly used character classes can be referred to by name (alpha, lower, upper, alnum, digit, punct, cntrl)
- Syntax [: name :]
- ▶ [a zA Z] is equivalent [[: alpha :]]
- [a zA Z0 9] is equivalent [[: alnum :]]
- ▶ [45a z] is equivalent [45[: lower :]]
- Important for portability across languages

Anchor Characters

- Anchors are used to match at the beginning or end of a line (or both).
- means beginning of the line
- \$ means end of the line



(D) (0) (2) (2) (2)



101 (B) (2) (2) (2) (2)

Subexpressions

- If you want to group part of an expression so that * or { } applies to more than just the previous character, use () notation
- Subexpressions are treated like a single character
- a* matches 0 or more occurrences of a
- abc* matches ab, abc, abcc, abccc, ...
- (abc)* matches abc, abcabc, abcabcabc, ...
- (abc)2,3 matches abcabc or abcabcabc

Global Regular Expressions Print – grep

- grep comes from the ed (Unix text editor) search command "global regular expression print" or g/re/p
- This was such a useful command that it was written as a standalone utility
- There are two other variants, egrep and fgrep that comprise the grep family
- grep is the answer to the moments where you know you want the file that contains a specific phrase but you can't remember its name

・ロン ・伊 とくさい くたい

Syntax

- Regular expression concepts we have seen so far are common to grep
- ▶ grep: \(and \), \{ and \}



・ロン ・西 ・ ・ さ ・ ・ き ・ ・ き