

# Hepatitis B Virus (HBV) amino acid alignments and genotype-specific consensus sequences

Philippa C. Matthews

Wellcome Trust Clinical Research Fellow, Nuffield Department of Medicine, University of Oxford,

Peter Medawar Building for Pathogen Research, South Parks Road, Oxford OX1 3SY, UK

Honorary Consultant in Microbiology and Infectious Diseases, Oxford University Hospitals NHS Foundation Trust, John Radcliffe Hospital, Headington, Oxford OX3 9DU, UK

Email: [p.matthews@doctors.org.uk](mailto:p.matthews@doctors.org.uk)

## Background:

Hepatitis B Virus (HBV) is the prototype human hepadnavirus. Unified reference sequences are important for informing phylogenetic analysis, studying diversity within and between hosts, identifying sequence changes (sites of polymorphisms / insertions / deletions) and providing standardised sequence numbering. We have developed a database of HLA Class I epitopes within HBV, 'hepitopes', (on-line at <http://www.expmedndm.ox.ac.uk/hepitopes>), for which we required a unified approach to HBV sequence numbering in order to provide a reference by which epitopes can be located. This work provides aligned and numbered HBV sequences, subdivided by genotype where possible, in the form of a pdf file and an xls spreadsheet.

## Methods:

- Recently published HBV reference sequences<sup>1</sup> were used as a baseline. The sequences were retrieved from Genbank (<https://www.ncbi.nlm.nih.gov/>) where they can be found under the following accession numbers:
  - Pol: AJW31599
  - Large HBs: AJW31600
  - HBx: AJW31597
  - Pre-core: AJW31591 (first 29 residues only, as remainder of sequence is Core protein)
  - Core: AJW31598
- Amino acid sequences were downloaded from The Hepatitis B Virus Database (<https://hbvdb.ibcp.fr>) for large HBs, HBx, Core and Pol, for all proteins for which >100 sequences were available. (Download dates 30-Sept to 18-Oct 2016; see Table for numbers of sequences).
- Downloaded sequences were converted to fasta format using [https://hcv.lanl.gov/cgi-bin/FORMAT\\_CONVERSION/convert.cgi](https://hcv.lanl.gov/cgi-bin/FORMAT_CONVERSION/convert.cgi)
- Sequences were aligned using clustal omega (up to a maximum of 2000 sequences per alignment) at <http://www.ebi.ac.uk/Tools/msa/clustalo/>
- A consensus sequence and entropy data was generated from each alignment using tools available on-line at the Los Alamos website [https://www.hiv.lanl.gov/content/sequence/ENTROPY/entropy\\_one.html](https://www.hiv.lanl.gov/content/sequence/ENTROPY/entropy_one.html)
- Alignments for sequences for Genotype E Pol and the first 29 residues of genotype D pre-core protein failed to align, so these are not represented in the consensus sequences.

Protein	Geno A	Geno B	Geno C	Geno D	Geno E	Geno F
Pol	1020	1403	1869	982	272	229
Large surface	1271	1799	1897	1167	292	237
HBx	862	1224	2000	1058	288	340
Pre-core	984	1149	1738	768	313	219
Core	1371	1500	1991	1433	453	252

**Table:** Number of sequences informing consensus for each HBV protein / genotype

## References:

- Liu WC, Lin CP, Cheng CP, Ho CH, Lan KL, Cheng JH, et al. Aligning to the sample-specific reference sequence to optimize the accuracy of next-generation sequencing analysis for hepatitis B virus. *Hepatol Int*. 2016 Jan;10(1):147-57.

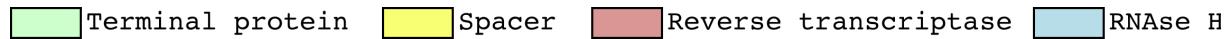
**HBV Pol Amino Acid Alignment (page 1/5)**Reference sequence is from Liu WC, et al., *Hepatol. Int.* 2016;10:147-57.Consensus sequences are derived from all available sequences at (<https://hbvdb.ibcp.fr>); see methods (page 1).

'-' designates a residue the same as that within the reference sequence.

'.' designates missing data

Residues that differ from the reference sequence are specified.

The sequence is numbered consecutively from 1 and divided into blocks of ten residues (dashed lines) for ease of navigation.


**Position** 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50

<b>Reference</b>	M P L S Y Q H F R K   L L L D E E A G P   L E E E L P R L A D   E G L N R R V A E D   L N L G N P N V S I
<b>Geno A</b>	- - - - - D - - - - -   - - - - -   A D - - - - -   L - - - - -
<b>Geno B</b>	- - - - - - - - - - -   - - - - -   - - - - - - - - - - -   L - - - - -
<b>Geno C</b>	- - - - - D - - - - -   - - - - -   - - - - - - - - - - -   L - - - - -
<b>Geno D</b>	- - - - - R - - - - -   D - - - - -   - - - - - - - - - - -   L - - - - -
<b>Geno F</b>	- - - - - P - - - - -   D - - - - -   - - - - - - - - - - -   Q L - - - - -

**Position** 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100

<b>Reference</b>	P W T H K V G N F T   G L Y S S T V P C F   N P K W Q T P S F P D I H L Q E D I V D R C E Q F V G P L T
<b>Geno A</b>	- - - - - I - - - E - - - - -   K - - - - -   I N - - Q - - - - -
<b>Geno B</b>	- - - - - - - - - - -   - - - - -   - - - - - - - - - - -   K - - - - -
<b>Geno C</b>	- - - - - V - - E - - - H - - - I N - - Q - Y - - - - -
<b>Geno D</b>	- - - - - V - - H - K - - N - - H Q - - I K K - - - - -
<b>Geno F</b>	- - - - - A - - D - L - - - - - H Q - L I S K - - - - -

**Position** 101 102 103 104 105 106 107 108 109 110 111 112 113 114 115 116 117 118 119 120 121 122 123 124 125 126 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143 144 145 146 147 148 149 150

<b>Reference</b>	V N E T R R L K L I   M P A R F Y P N V T   K Y L P L D K G I K   P Y Y P E H V V N H   Y F Q T R H Y L H T
<b>Geno A</b>	- - - K - - - - - T H - - - - -   - - - D Q - - - - -
<b>Geno B</b>	- - - N -
<b>Geno C</b>	- - - K - - - - L - - - - - A - - K - - - - -
<b>Geno D</b>	- - - K - - Q - - - - - - - - - - - L - - - - -
<b>Geno F</b>	K - - L - - - V - - - F - K - - - F - M E - - - - - K - - - - -

**Position** 151 152 153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168 169 170 171 172 173 174 175 176 177 178 179 180 181 182 183 184 185 186 187 188 189 190 191 192 193 194 195 196 197 198 199 200

<b>Reference</b>	L W K A G I L Y K R   E S T R S A S F C G   S P Y S W E Q D L Q   H G R L V F Q T S K   R H G D K S F C P Q
<b>Geno A</b>	- - - - - T - - - - -   E - - - - -   I K - - Q - - - E - - S -
<b>Geno B</b>	- - - - - - - - - - -   - - - - -   - - - - - - - - - - -
<b>Geno C</b>	- - - T - - - - E - - - - - Q - T - - E - - S -
<b>Geno D</b>	- - - - - T - H - - - - E - - - - - - - - - - - A E - - H Q -
<b>Geno F</b>	- - - - - - - - - - - E - - - S T S L N D - - G - - T E - L - A -

 continues on next page

**HBV Pol Amino Acid Alignment (page 2/5)**Reference sequence is from Liu WC, et al., *Hepatol. Int.* 2016;10:147-57.Consensus sequences are derived from all available sequences at (<https://hbvdb.ibcp.fr>); see methods (page 1).

'-' designates a residue the same as that within the reference sequence.

'.' designates missing data

Residues that differ from the reference sequence are specified.

The sequence is numbered consecutively from 1 and divided into blocks of ten residues (dashed lines) for ease of navigation.


**Position**

201 202 203 204 205 206 207 208 209 210 211 212 213 214 215 216 217 218 219 220 221 222 223 224 225 226 227 228 229 230 231 232 233 234 235 236 237 238 239 240 241 242 243 244 245 246 247 248 249 250 251

Reference	S P G I L R G S S V	G P S I Q S Q L R K	S R L G P Q P A Q G	Q L A G R Q Q G G S	G S I R A R V H P S
Geno A	P S - - - S R - - -	C - R - - - K Q - - -	L - H - - - P - - - S S Q P - R - - -	A - - - - - - - - - - -	
Geno B	- - - - - P R - - -	C - - - - - - - - - - -	- - - - - - - - - - -	R - - - - - - - - - - -	
Geno C	- S - - - S R P - -	C V R - - - K Q - - -	L - Q - - - S - - - R G K S - R - - -	T - - - - - - - - - - -	
Geno D	- S - - - S R P P - -	S - L - - K H - - -	L - S Q - - - H - R - - - R - W - - -	G I - T - - - - - - - - -	
Geno F	- S - - - S R P - A -	S - - - G K F Q Q - - -	L - Q K - - - N G K - - - R L - S - - - T P	- - - - - - - - - - - - -	

**Position**

251 252 253 254 255 256 257 258 259 260 261 262 263 264 265 266 267 268 269 270 271 272 273 274 275 276 277 278 279 280 281 282 283 284 285 286 287 288 289 290 291 292 293 294 295 296 297 298 299 300

Reference	P W G T V G V E P S G	S G Q T H N C A S	S S S C L H Q S A V R K A A Y S L I S T	T S K G H S S S G H
Geno A	T R R Y F - - - - -	H I D H S V N N - - - - -	- - - - - H L - - - R Q - - -	- - - - - - - - - - -
Geno B	- - - - - - - - - - -	P T - - - - - - - - - - -	- - - - - - - - - - -	- - - - - - - - - - -
Geno C	T R R S F - - - - -	H I D - S - - T - - -	- - - T - - H L - - - R Q - - -	- - - - - - - - - - -
Geno D	A R R P F - - - - -	H T T - L - - K - A - - Y - P - - -	- - - P A V - - F E K - - -	- - - - - - - - - - -
Geno F	T R W P - - - - -	T - C A N - L - - R - A - - F - - -	- - - E K - N P S L - - - R - T - T - -	- - - - - - - - - - -

**Position**

301 302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350

Reference	A V E L H H F P P N	S S R S Q S P G P V	P S C W W L Q F R N	S E P C S E Y C L C H I V N L I D D W G
Geno A	- - - F - C L - - S -	A G - - - Q - S - F - - -	- - - K - - - - - - -	S - L - - - R - - -
Geno B	- - - - - - - - - - -	Q - - - L - - - - - -	- - - - - - - - - - -	- - - - - - - - - - -
Geno C	- - - - N I - - S - A - - - E - I L - - -	K - - - D - - - T - - -	L - - - - - - - - - - -	
Geno D	- - - - N L - - - A - - - E R - - -	F P - - - K - - - D - - -	S - - - - - - - - - - -	
Geno F	- - - - N S V - - G - V - - E G K - S - F - - -	D T - - - D - - - S - - - I - - L - - -	R - - - - - - - - - - -	

**Position**

351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 378 379 380 381 382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400

Reference	P C T E H G E H L I	R T P R T P A R V T	G G V F L V D K N P	H N T T E S R L V V D F S Q F S R G N T
Geno A	- - D - - - H - I - - -	- - - - - - - - - - -	- - A - - - - - - -	I - - - - - - - - - - -
Geno B	- - - - - R - - - - -	- - - - - - - - - - -	- - - - - - - - - - -	- - - - - - - - - - -
Geno C	- - - - - N - I - - -	- - - - - - - - - - -	- - - - - - - - - - -	S - - - - - - - - - - -
Geno D	- - A - - - H - I - - -	- - - - - - - - - - -	- - A - - - - - - -	Y - - - - - - - - - - -
Geno F	- - Y - - - Q - H - - -	- - - - - - - - - - -	- - - - - - - - - - -	T - - - - - - - - - - -

→ continues on next page

## HBV Pol Amino Acid Alignment (page 3/5)

Reference sequence is from Liu WC, et al., *Hepatol. Int.* 2016;10:147-57.

Consensus sequences are derived from all available sequences at (<https://hbvdb.ibcp.fr>); see methods (page 1).

'-' designates a residue the same as that within the reference sequence.

‘.’ designates missing data

Residues that differ from the reference sequence are specified.

The sequence is numbered consecutively from 1 and divided into blocks of ten residues (dashed lines) for ease of navigation. YMDD motif is marked with a box (residues 549-552 in this alignment).

Terminal protein      Spacer      Reverse transcriptase      RNase H

<b>Position</b>	501 502 503 504 505 506 507 508 509 510 511 512 513 514 515 516 517 518 519 520 521 522 523 524 525 526 527 528 529 530 531 532 533 534 535 536 537 538 539 540 541 542 543 544 545 546 547 548 549 550
<b>Reference</b>	L H L Y S H P I I L G F R K I P M G V G L S P F L L A Q F T S A I C S V V R R A F F P H C L A F S Y M
<b>Geno A</b>	- - - - - V - - - - -
<b>Geno B</b>	- - - - - - - - - - -
<b>Geno C</b>	- - - - - - - - - - -
<b>Geno D</b>	- - - - - - - - - - -
<b>Geno F</b>	- - - - - - - - - - -

<b>Position</b>	551 552 553 554 555 556 557 558 559 560 561 562 563 564 565 566 567 568 569 570 571 572 573 574 575 576 577 578 579 580 581 582 583 584 585 586 587 588 589 590 591 592 593 594 595 596 597 598 599 600	D D V V L G A K S V Q H L E S L Y A A V T N F L L S L G I H L N P H K T K R W G Y S L N F M G Y V I
<b>Reference</b>		R - T - L - N - I
<b>Geno A</b>		R - T - L - N - I
<b>Geno B</b>		L
<b>Geno C</b>		F T S I L N
<b>Geno D</b>		F T L N H
<b>Geno F</b>	L T	T S T H

→ continues on next page

## HBV Pol Amino Acid Alignment (page 4/5)

Reference sequence is from Liu WC, et al., *Hepatol. Int.* 2016;10:147-57.

Consensus sequences are derived from all available sequences at (<https://hbvdb.ibcp.fr>); see methods (page 1).

'-' designates a residue the same as that within the reference sequence.

‘.’ designates missing data

Residues that differ from the reference sequence are specified.

The sequence is numbered consecutively from 1 and divided into blocks of ten residues (dashed lines) for ease of navigation.



→ continues on next page

HBV Pol Amino Acid Alignment (page 5/5)

Reference sequence is from Liu WC, et al., *Hepatol. Int.* 2016;10:147-57.

Consensus sequences are derived from all available sequences at (<https://hbvdb.ibcp.fr>); see methods (page 1).

'-' designates a residue the same as that within the reference sequence.

‘.’ designates missing data

Residues that differ from the reference sequence are specified.

The sequence is numbered consecutively from 1 and divided into blocks of ten residues (dashed lines) for ease of navigation.

	Terminal protein	Spacer	Reverse transcriptase	RNAse H	
Position	801 802 803 804 805 806 807 808 809	810 811 812 813 814 815 816 817 818 819	820 821 822 823 824 825 826 827 828 829	830 831 832 833 834 835 836 837 838 839	840 841 842 843
Reference	L L R L L Y R P T T G R T S L Y A D S P S V P S H L P D R V H F A S P L H V A W K P P				
Geno A	--P F Q--	-V-	-V-		
Geno B	--				
Geno C	--H-P F--	-V-			
Geno D	--P F--				
Geno F	--P F Q--				

## **HBV Large S Protein Amino Acid Alignment (page 1/2)**

Reference sequence is from Liu WC, et al., *Hepatol. Int.* 2016;10:147-57.

Consensus sequences are derived from all available sequences at (<https://hbvdb.ibcp.fr>); see methods (page 1).

'-' designates a residue the same as that within the reference sequence.

‘.’ designates missing data

Residues that differ from the reference sequence are specified.

The sequence is numbered consecutively from 1 and divided into blocks of ten residues (dashed lines) for ease of navigation.

## **HBV Large S Protein Amino Acid Alignment (page 2/2)**

Reference sequence is from Liu WC, et al., *Hepatol. Int.* 2016;10:147-57.

Consensus sequences are derived from all available sequences at (<https://hbvdb.ibcp.fr>); see methods (page 1).

'-' designates a residue the same as that within the reference sequence.

‘.’ designates missing data

Residues that differ from the reference sequence are specified.

The sequence is numbered consecutively from 1 and divided into blocks of ten residues (dashed lines) for ease of navigation.

<b>Position</b>	351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 378 379 380 381 382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400
<b>Reference</b>	V P F V Q W F V G L S P T V W L S V I W M M W Y W G P S L Y N I L S P F M P L L P I F F C L W V Y I
<b>Geno A</b>	- - - - - A - - - - - S - V - - I - - - - -
<b>Geno B</b>	- - - - - - - - - F - - - - - - - - -
<b>Geno C</b>	- - - - - - - - - - - - - - L - - - - -
<b>Geno D</b>	- - - - - - - - - S - - L - - - - -
<b>Geno E</b>	- - - - - A - - - - - - - - I - - - - -
<b>Geno F</b>	- Q - - - C - - - - L - - I - - - C S - - I - - - - C Y - - S -

## **HBV X Protein Amino Acid Alignment**

Reference sequence is from Liu WC, et al., *Hepatol. Int.* 2016;10:147-57.

Consensus sequences are derived from all available sequences at (<https://hbvdb.ibcp.fr>); see methods (page 1).

'-' designates a residue the same as that within the reference sequence.

‘.’ designates missing data

Residues that differ from the reference sequence are specified.

The sequence is numbered consecutively from 1 and divided into blocks of ten residues (dashed lines) for ease of navigation.

<b>Position</b>	151
<b>Reference</b>	F T S A
<b>Geno A</b>	- - - -
<b>Geno B</b>	- - - -
<b>Geno C</b>	- - - -
<b>Geno D</b>	- - - -
<b>Geno E</b>	- - - -
<b>Geno F</b>	- - - -

## HBV Pre-Core Protein Amino Acid Alignment (first 29 residues only)

The remaining residues (30 onwards) are represented in the alignment for core sequence (see page 11)

Reference sequence is from Liu WC, et al., *Hepatol. Int.* 2016;10:147-57.

Consensus sequences are derived from all available sequences at (<https://hbvdb.ibcp.fr>); see methods (page 1).

'-' designates a residue the same as that within the reference sequence.

',' designates missing data

Residues that differ from the reference sequence are specified.

The sequence is numbered consecutively from 1 and divided into blocks of ten residues (dashed lines) for ease of navigation.

Position	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29
Reference	M	Q	L	F	H	L	C	L	I	I	S	C	S	C	P	T	V	Q	A	S	K	L	C	L	G	W	L	W	G
Geno A	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
Geno B	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
Geno C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
Geno D	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
Geno F	-	-	-	-	-	-	-	-	-	-	F	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	

## HBV Core Protein Amino Acid Alignment

The pre-core protein incorporates an additional 29 residues upstream of this sequence (see sequences represented on page 10).

Reference sequence is from Liu WC, et al., *Hepatol. Int.* 2016;10:147-57.

Consensus sequences are derived from all available sequences at (<https://hbvdb.ibcp.fr>); see methods (page 1).

'-' designates a residue the same as that within the reference sequence.

‘.’ designates missing data

Residues that differ from the reference sequence are specified.

The sequence is numbered consecutively from 1 and divided into blocks of ten residues (dashed lines) for ease of navigation.

Position	151	152	153	154	155	156	157	158	159	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175	176	177	178	179	180	181	182	183
Reference	R	R	G	R	S	P	R	R	R	T	P	S	P	R	R	R	S	Q	S	P	R	R	R	S	Q	S	R	E	S	K	C		
Geno A	--	D	.	R	G	-	.	.	.	-	-	-	-	T	P	.	-	-	-	-	-	-	-	-	-	-	-	-	-	-			
Geno B	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-			
Geno C	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-			
Geno D	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-			
Geno E	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	P	A	-	-		
Geno F	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	P	A	-	-