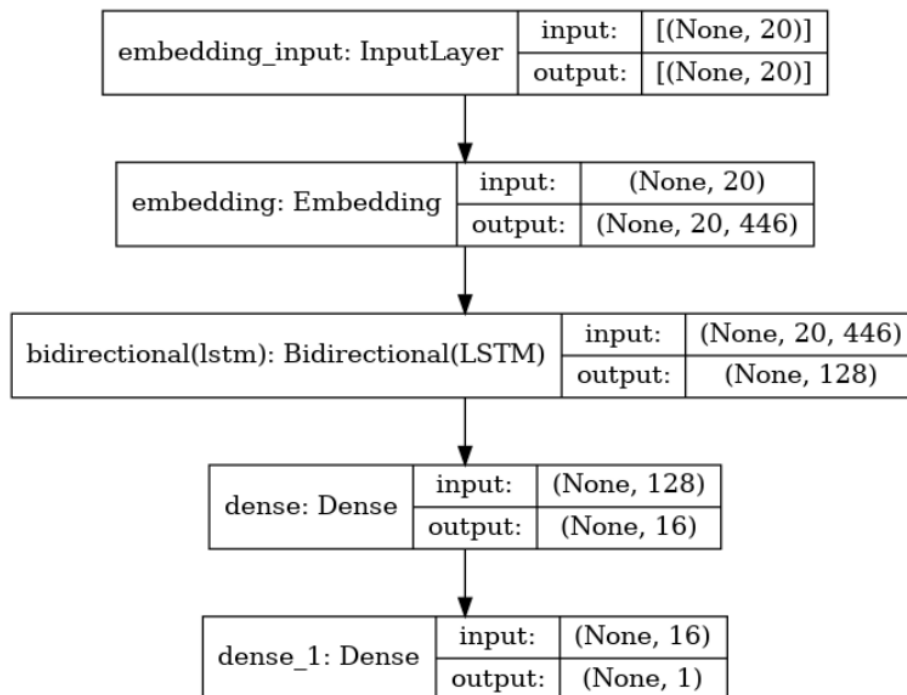




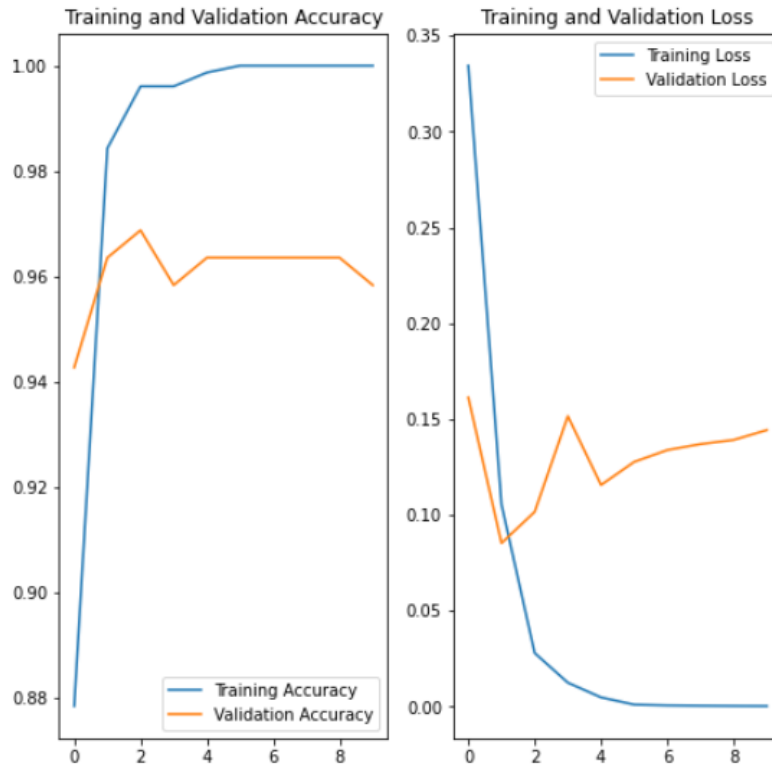
Dari hasil ini, terlihat kata “*free*”, “*mobile*”, “*call*”, dan “*claim*” menjadi kata yang paling banyak muncul dalam kelas spam. Hal ini menunjukkan data dapat digunakan untuk membentuk model pembelajaran mesin mengingat kata-kata tersebut sering muncul di dalam email spam.

Selanjutnya, dataset pelatihan dibagi kembali menjadi dataset pelatihan sejumlah 765 dan sisanya masuk ke dalam dataset validasi. Tokenisasi dan *padding* diterapkan pada dataset untuk meningkatkan performa model. Tokenisasi adalah tahap pembagian kalimat menjadi bagian-bagian kecil. Bagian-bagian kecil tersebut kemudian dikodekan menjadi suatu angka tertentu yang dapat diproses secara matematik.

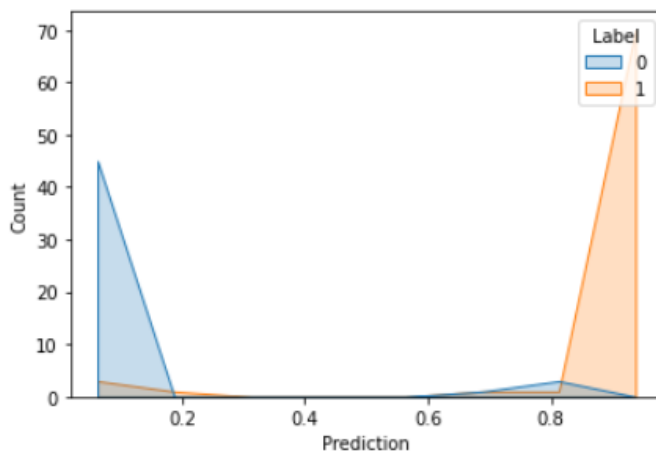
Fungsi loss yang digunakan adalah *binary cross entropy* dengan algoritma pengoptimal berupa *adam* dan metrik yang digunakan adalah akurasi. Berikut arsitektur model yang digunakan.



Hasil dari pelatihan dapat dilihat pada gambar dibawah.

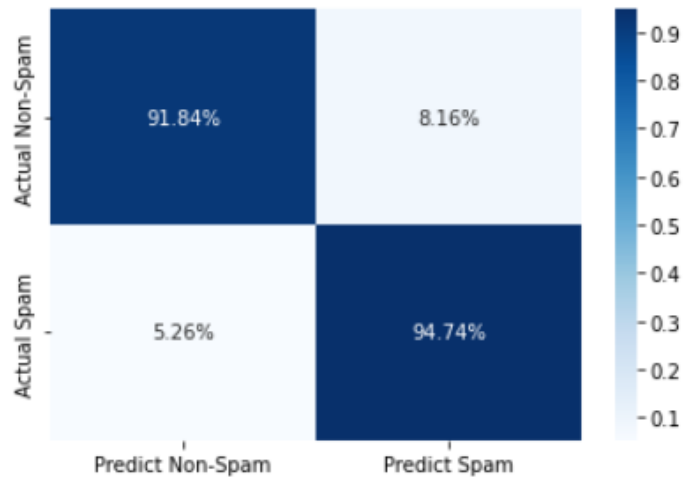


Performa model dapat mencapai 100% akurasi pada dataset training dan 96% akurasi pada dataset validasi. Model ini kemudian digunakan untuk memprediksi email pada dataset tes yang belum dilihat oleh model.



Visualisasi di atas menunjukkan distribusi hasil klasifikasi model LSTM yang telah dibuat terhadap dataset tes. Dapat dilihat model berhasil mengklasifikasikan email spam, yang

ditandai dengan label 1, dan email tidak spam. Nilai *precision*, *recall*, dan *f1 score* sejumlah 0.93, 0.93, dan 0.93.



Nilai matriks *confusion* dari model yang telah dibuat ditunjukkan oleh gambar di atas. Model berhasil memprediksi 94.74% email spam dan 91.84% email bukan spam.

Kesimpulan

Telah dibuat model pembelajaran mesin LSTM untuk mengklasifikasikan email spam. Nilai *precision*, *recall*, dan *f1 score* sejumlah 0.93, 0.93, dan 0.93. Model berhasil memprediksi 94.74% email spam dan 91.84% email bukan spam.